

Joint Representation Classification for Collective Face Recognition*

Liping Wang[†] Songcan Chen[‡]

October 14, 2018

Abstract

Sparse representation based classification (SRC) is popularly used in many applications such as face recognition, and implemented in two steps: representation coding and classification. For a given set of testing images, SRC codes every image over the base images as a sparse representation then classifies it to the class with the least representation error. This scheme utilizes an individual representation rather than the collective one to classify such a set of images, doing so obviously ignores the correlation among the given images. In this paper, a joint representation classification (JRC) for collective face recognition is proposed. JRC takes the correlation of multiple images as well as a single representation into account. Under the assumption that the given face images are generally related to each other, JRC codes all the testing images over the base images simultaneously to facilitate recognition. To this end, the testing inputs are aligned into a matrix and the joint representation coding is formulated to a generalized $l_{2,q} - l_{2,p}$ -minimization problem. To uniformly solve the induced optimization problems for any $q \in [1, 2]$ and $p \in (0, 2]$, an iterative quadratic method (IQM) is developed. IQM is proved to be a strict descent algorithm with convergence to the optimal solution. Moreover, a more practical IQM is proposed for large-scale case. Experimental results on three public databases show that the JRC with practical IQM not only saves much computational cost but also achieves better performance in collective face recognition than the state-of-the-arts.

Keywords: SRC; JRC; IQM; practical IQM.

1 Introduction

Recently, representation coding based classification and its variants have been developed for face image recognition (FR) [1–5]. This schemes achieve a great success in FR

*The work is partially supported by the Chinese grants NSFC11471159, NSFC61170151 and Natural Science Foundation of Jiangsu Province (BK20141409).

[†]Department of Mathematics, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China. Email: wlpmath@nuaa.edu.cn.

[‡]Department of Computer Science and Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing, 210016, China. Email: s.chen@nuaa.edu.cn.

and boost the applications of image classification [6, 7]. Sparse representation based classification (SRC) [1] is the most known one which directly uses the sparse code for classification and efficiently recognizes the class giving the most compact representation. The main idea can be summarized to two steps: 1) coding a testing sample as a linear combination of all the training samples, then 2) classifying the testing sample to the most compact one by evaluating coding errors. Typical SRC employs the following l_1 -minimization as the sparse representation model,

$$\min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \varepsilon, \quad (1)$$

where $A \in R^{m \times d}$ is the dictionary of coding atoms and $y \in R^m$ is a given observation. $x \in R^d$ is the coding vector and $\varepsilon > 0$ denotes a noisy level. SRC outputs the identity of y as

$$\text{identity}(y) = \arg \min_{1 \leq i \leq I} \{\|y - Ax_i^*\|_2\}, \quad (2)$$

where I denotes the number of classes and x_i^* is the coding coefficient vector associated with class i . The experimental results reported in [1] exhibit that SRC scheme achieves amazing performance. But the authors of [2] argued that SRC over emphasized the importance of l_1 -norm sparsity but ignored the effect of collaborative representation. Consequently, a collaborative representation based classification with regularized least square (CRC-RLS) was presented in [2] for face recognition

$$\min_x \|x\|_2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \varepsilon. \quad (3)$$

Anyway, problem (3) is easier to solve than (2) for its smoothness. Models (2) and (3) can be considered as the least square problems with different regularizers,

$$\min_x \|y - Ax\|_2^2 + \lambda \|x\|_1 \quad \text{and} \quad \min_x \|y - Ax\|_2^2 + \lambda \|x\|_2^2. \quad (4)$$

Moreover, Wright et al. [3] ever used variant l_1 -norm to improve the coding fidelity of y over A ,

$$\min_x \|y - Ax\|_1 + \lambda \|x\|_1. \quad (5)$$

Actually, the models (2)-(5) can be uniformly included in the framework

$$\min_x \|y - Ax\|_q^q + \lambda \|x\|_p^p, \quad 1 \leq q \leq 2, \quad 0 < p \leq 2. \quad (6)$$

In (6), the representation and regularization measurements are extended to be $\|\cdot\|_q$ ($1 \leq q \leq 2$) and $\|\cdot\|_p$ ($0 < p \leq 1$) respectively. This modification provides possibility to adaptively choose the most suitable model for different applications. Moreover, the computational experiences [13–15] have showed that fractional norm l_p ($0 < p < 1$) exhibits sparser pattern than l_1 -norm. The unified generalization formula (6) is expected to achieve better performance. On the other hand, model (6) is a vector representation based framework which implies the following weaknesses.

- Model (6) uses coding vector to represent testing samples one by one. In many face recognition, a great of number of images for each known subject have been collected from video sequence or photo album. The face recognition has to be conducted

with a set of probe images rather than a single one [8]. In this case, representation coding based classification like model (6) can not efficiently work.

- Any testing sample is coded independently from each other in (6). This approach takes no account of the correlation hidden in the image set. The difference and similarity between multiple pictures are totally ignored. It is well known that the collective faces share some similar feature patterns, such as eye or mouth pixels is more powerful in discrimination than those of forehead or cheek.

- When q, p in (6) take different values, the involved optimization problems have to be solved by different algorithms. For example, (1) is solved by $l_1 - l_s$ solver [9] or alternative direction of multiplier method while (3) chooses the algorithm presented in [2].

To overcome the weaknesses in (6) and make sufficient use of collective relationship among the given set of images, we consider to jointly represent all the test samples simultaneously over the training sample base. Here we employ matrix instead of vector as the coding variable to evaluate the distribution of feature space. This idea induces a joint representation based classification (JRC) for collective face recognition and reduces it to a $l_{2,q} - l_{2,p}$ -minimization. To solve the derived optimization problem, a unified algorithm is designed and its convergence behavior is also analyzed. Experiments on three public face datasets validate the improvement of JRC over the state-of-the-arts.

This paper is organized as follows. In the second section, a joint representation based classification (JRC) will be established. The third section is dedicated to a unified algorithm for solving the special optimization problem induced by JRC. Some computational details are considered in the fourth section and an improved practical algorithm is proposed. The experimental results are reported in the last section.

2 Joint Representation Classification for Collective Face Recognition

2.1 Joint Representation Model

Suppose that we have I classes of subjects in the dataset. $A_i \in R^{m \times d_i} (1 \leq i \leq I)$ denotes the i -th class, and each column of A_i is a sample of class i . Hence all the training samples are aligned by $A = [A_1, A_2, \dots, A_I] \in R^{m \times d}$, where $d = \sum_{i=1}^I d_i$. Given a collection of query images $y_1, y_2, \dots, y_n \in R^m$, model (6) codes each $y_j (1 \leq j \leq n)$ by the training samples A as

$$y_j \approx Ax_j, \tag{7}$$

where $x_j \in R^d$ is the coding vector associated with y_j . If y_j is from the i -th class, then A_i is the most compact representation dictionary and the optimal solution x_j^* to (6) can be used for classification. Obviously, coding pattern (7) depends on the single test sample y_j individually for classification but takes no account of the correlation with other samples ($y_l, l \neq j$). Even though different frontal faces take on different

appearances, they share similar features such as two eyes and brows at the upper face while nose and mouth at the lower. Difference and similarity of multiple face pictures form a unitary feature of the given set of images which play an important role for collective face recognition.

Denote $Y = [y_1, y_2, \dots, y_n] \in R^{m \times n}$ all the query images, we propose to jointly represent the image set simultaneously by

$$Y \approx AX, \quad (8)$$

where $X = [x_1, x_2, \dots, x_n] \in R^{d \times n}$ stands for the collective coding matrix. As far as the columns are concerned, system (8) is an easy consequence of (7). To measure the fidelity of the joint coding system (8), we consider X in another sense. Let $A^i \in R^d$ and $Y^i \in R^n$ be the i -th ($i = 1, 2, \dots, m$) row vectors of matrix A and Y respectively, formula (8) is equivalent to

$$X^T(A^i)^T \approx (Y^i)^T \quad \text{for } i = 1, 2, \dots, m. \quad (9)$$

It is noticed that A and Y array the sampled images column by column, hence their rows span the feature space. In feature extraction view, the collective coding matrix X also projects the training feature space to approximate the testing feature space. Traditional least square regression aims to minimize the error

$$\min_X \sum_{i=1}^m \|X^T(A^i)^T - (Y^i)^T\|_2^2 \quad \text{or} \quad \min_X \sum_{i=1}^m \|A^i X - Y^i\|_2^2. \quad (10)$$

Actually (10) can be easily reformulated as

$$\min_X \sum_{i=1}^m \|(AX - Y)^i\|_2^2, \quad (11)$$

where $(AX - Y)^i$ is the i -th row vector of $AX - Y$. Especially when the number of column in $AX - Y$ is 1, the formula (11) is reduced to the fidelity function of (4). Then we prefer a uniform generalization of (4) and (5) in the sense

$$\sum_{i=1}^m \|(AX - Y)^i\|_2^q, \quad (1 \leq q \leq 2). \quad (12)$$

Under the assumption that joint representation and feature distribution share the similar pattern for all testing face images, we use the following regularization

$$\sum_{i=1}^d \|X^i\|_2^p, \quad (0 < p \leq 2), \quad (13)$$

where X^i is the i -th row vector of X for $i = 1, 2, \dots, d$. Combining (12) and (13), we present the joint representation model for classification as follows

$$\min_X \sum_{i=1}^m \|(AX - Y)^i\|_2^q + \lambda \sum_{i=1}^d \|X^i\|_2^p, \quad (1 \leq q \leq 2, 0 < p \leq 2). \quad (14)$$

When the number of column in Y is 1, model (14) is reduced to coding vector version (6). Compared with coding vector x , joint coding matrix X unites sample representation with feature projection which somewhat reflects the integral structure of dataset. Hence (14) is a general extension of (3)-(6). To simplify the formulation, we introduce the mixed matrix norm $l_{2,p}$ ($p > 0$) (taking $\|X\|_{2,p}$ for example)

$$\|X\|_{2,p} = \left(\sum_{i=1}^d \|X^i\|_2^p \right)^{\frac{1}{p}}, \quad X \in R^{d \times n}, \quad (15)$$

where X^i denotes the i -th row of X . Then (14) is rewritten as

$$\min_X \|AX - Y\|_{2,q}^q + \lambda \|X\|_{2,p}^p, \quad (1 \leq q \leq 2, 0 < p \leq 2). \quad (16)$$

Especially when $p \in (0, 1)$, $l_{2,p}$ is not a valid matrix norm because it does not satisfy the triangular inequality of matrix norm axioms. Meanwhile the involved fractional matrix norm based minimization (16) is neither convex nor Lipschitz continuous which brings computational challenge. Designing an efficient algorithm for such $l_{2,q} - l_{2,p}$ -minimizations is very important. It is also the most challenging task in this paper.

2.2 Joint Representation Based Classification

For fixed parameter q and p , suppose that X^* is a minimizer of optimization problem (16), that is

$$X^* = \arg \min_X \|AX - Y\|_{2,q}^q + \lambda \|X\|_{2,p}^p. \quad (17)$$

If X^* is partitioned to I blocks as follows

$$X^* = \begin{bmatrix} X_1^* \\ \vdots \\ X_i^* \\ \vdots \\ X_I^* \end{bmatrix}, \quad (18)$$

where $X_i^* \in R^{d_i \times n}$ ($1 \leq i \leq I$). Let \hat{X}_i^* denote the coding matrix associated with class i , that is

$$\hat{X}_i^* = \begin{bmatrix} 0 \\ \vdots \\ X_i^* \\ \vdots \\ 0 \end{bmatrix}, \quad (19)$$

then $A\hat{X}_i^* = A_i X_i^*$ ($1 \leq i \leq I$). For each testing image y_j ($j = 1, 2, \dots, n$), we classify y_j to the class with the most compact representation. By evaluating the error corresponding to each class

$$\|(Y - A\hat{X}_i^*)_j\|_2, \quad i = 1, 2, \dots, I \quad (20)$$

we pick out the index outputting the least error. The joint representation based classification for face recognition can be concluded as follows.

Algorithm 2.1. (*JRC scheme for FR*)

1. *Start:* Given $A \in R^{m \times d}$, $Y \in R^{m \times n}$ and select parameters $\lambda > 0$, $q \in [1, 2]$ and $p \in (0, 2]$.
2. Solve $l_{2,q} - l_{2,p}$ -minimization problem (16) for coding matrix X^* .
3. *For* $j = 1 : n$
 For $i = 1 : I$
 $e_i(y_j) = \|(Y - A_i X_i^*)_j\|_2$
 end
 Identity $(y_j) = \arg \min_{1 \leq i \leq I} \{e_i(y_j)\}$
 end

When $n = 1$, observation Y contains only a single testing sample and JRC is reduced to vector representation based classification. Further on, SRC, CRC-RLS and l_1 -norm fidelity model (5) are the special cases of JRC when $q = 2$ & $p = 1$, $q = p = 2$ and $q = p = 1$ respectively. In short, the main contributions of JRC lie in:

1. JRC implements collective face representation simultaneously. This scheme is more economical and efficient in computational cost and CPU time. Moreover, JRC can handle image set based face recognition which broadens the applications of vector representation based classifications.
2. Joint coding technique fuses the difference of each testing sample representation and the similarity hidden in the feature space of multiple face images. For example, when $0 < p \leq 1$ all query image are jointly represented by the training samples with the similarly sparse feature distribution.
3. In the next section, a uniform algorithm will be developed to solve the optimization problem (16) for any $q \in [1, 2]$ and $p \in (0, 2)$. The algorithm is strict decreasing until it converges to the optimal solution to problem (16). To the best of our knowledge, it is an innovative approach to solve such a generalized $l_{2,q} - l_{2,p}$ -minimization.

It is worth to point out that the JRC scheme can be easily extended for the presence of pixel distortion, occlusion or high noise in test images. Modify (8) as

$$Y = AX + E, \quad (21)$$

where $E \in R^{m \times n}$ is an error matrix. The nonzero entries of E locate the corruption or occlusion in Y . Substitute $\hat{A} = [A, I] \in R^{m \times (d+m)}$ and $\hat{X} = \begin{bmatrix} X \\ E \end{bmatrix} \in R^{(d+m) \times n}$ for A and X respectively, a stable joint coding model can be formulated to

$$\min_{\hat{X}} \|\hat{A}\hat{X} - Y\|_{2,q}^q + \lambda \|\hat{X}\|_{2,p}^p, \quad (1 \leq q \leq 2, 0 < p \leq 2). \quad (22)$$

Once a solution $\hat{X}^* = \begin{bmatrix} X^* \\ E^* \end{bmatrix}$ to (22) is computed, setting $Y^* = Y - E^*$ recovers a clean image from corrupted subject. To identify the testing sample y_j , we slightly modify the error of y_j with each subject $e_i(y_j) = \|(Y - E^* - A_i X_i^*)_j\|_2$. Thus a robust JRC is an easy consequence of Algorithm 2.1. The corresponding algorithm and theoretical analysis can be similarly demonstrated. This paper will not concentrate on this subject.

3 An Iterative Quadratic Method for JRC

Obviously, efficiently solving optimization problem (16) plays the most important role in scheme 2.1. The mentioned models (1), (3) and (5) are special cases of (16), the algorithms used in [1–3] to solve those special problems can not be directly extended. Such generally mixed matrix norm based minimizations as (16) have been widely used in machine learning. Rakotomamonjy and his co-authors [10] proposed to use the mixed matrix norm $l_{q,p}$ ($1 \leq q < 2, 0 < p \leq 1$) in multi-kernel and multi-task learning. But the induced optimization problems in [10] have to be solved separately by different algorithms with respect to $p = 1$ and $0 < p < 1$. For grouped feature selection, Suvrit [11] addressed a fast projection technique onto $l_{1,p}$ -norm balls particularly for $p = 2, \infty$. But the derived method in [11] does not match model (16). Similar joint sparse representation has been used for robust multimodal biometrics recognition in [12]. The authors of [12] employed the traditional alternating direction method of multipliers to solve the involved optimization problem. Nie et al. [16] applied $l_{2,0+}$ -norm to semi-supervised robust dictionary learning, while the optimization algorithm has not displayed definite convergence analysis.

In this section, a unified method will be developed to solve the $l_{2,q}-l_{2,p}$ -minimization (16) for any $1 \leq q \leq 2$ and $0 < p \leq 2$. Especially when $p \in (0, 1)$, (16) is neither convex nor non-Lipschitz continuous which results in much computational difficulties. Motivated by the idea of algorithm in [17] for solving $l_{2,p}$ ($0 < p \leq 1$)-based minimization, we design an iteratively quadratic algorithm for such $l_{2,q}-l_{2,p}$ -minimization. Moreover, the convergence analysis will be uniformly demonstrated.

3.1 An Iteratively Quadratic Method

After simply transformation, the definition of $\|X\|_{2,p}^p$ (15) can be rewritten as

$$\|X\|_{2,p}^p = Tr(X^T H X), \quad (23)$$

where

$$H = \begin{cases} \text{diag}\left\{\frac{1}{\|X^1\|_2^{2-p}}, \frac{1}{\|X^2\|_2^{2-p}}, \dots, \frac{1}{\|X^d\|_2^{2-p}}\right\}, & p \in (0, 2); \\ I, & p = 2, \end{cases} \quad (24)$$

and $Tr(\cdot)$ stands for trace operation. If denote

$$G = \begin{cases} \text{diag}\left\{\frac{1}{\|(AX-Y)^1\|_2^{2-q}}, \frac{1}{\|(AX-Y)^2\|_2^{2-q}}, \dots, \frac{1}{\|(AX-Y)^m\|_2^{2-q}}\right\}, & q \in [1, 2); \\ I, & q = 2, \end{cases} \quad (25)$$

the objective function of (16) can be reformulated to

$$J(X) := Tr((AX - Y)^T G(AX - Y)) + \lambda Tr(X^T H X). \quad (26)$$

Hence the KKT point of unconstrained optimization problem (16) is also the stationary point of $J(X)$,

$$\frac{\partial J(X)}{\partial X} = qA^T G(AX - Y) + \lambda p H X = 0, \quad (27)$$

solving (16) is reduced to find the solution to equations (27). If $A^T G A + \lambda \frac{p}{q} H$ is invertible, equation (27) is equivalent to

$$X = (A^T G A + \lambda \frac{p}{q} H)^{-1} A^T G Y. \quad (28)$$

To find the iterative solution to system (28), let us consider a closely related optimization problem

$$\min_X \hat{J}(X) := Tr((AX - Y)^T G(AX - Y)) + \lambda \frac{p}{q} Tr(X^T H X). \quad (29)$$

$\hat{J}(X)$ is almost equivalent to $J(X)$ in spite of a scaled factor $\frac{p}{q}$ in regularization parameter. If an iterative approximate solution X_k to (29) has been generated, G_k and H_k can be derived from X_k as definitions (24, 25). Then we can compute the next iterative matrix X_{k+1} by solving the following subproblem

$$\min_X Tr((AX - Y)^T G_k(AX - Y)) + \lambda \frac{p}{q} Tr(X^T H_k X). \quad (30)$$

Actually, (30) is a scaled quadratic approximation to $J(X)$ at the iterative point X_k . Let $M_k = A^T G_k A + \lambda \frac{p}{q} H_k$, since G_k and H_k are usually symmetric and positive definite, problem (30) is equivalent to the following quadratic optimization problem

$$\min_X Q_k(X) := \frac{1}{2} Tr(X^T M_k X) - Tr(Y^T G_k A X). \quad (31)$$

The minimizer to $Q_k(X)$ is also the solution to the linear system

$$M_k X = A^T G_k Y. \quad (32)$$

Based on the analysis and equations (23-32), the mixed $l_{2,q} - l_{2,p}$ ($1 \leq q \leq 2, 0 < p \leq 2$) norm based optimization problem (16) can be iteratively solved by a sequence of quadratic approximate subproblems. Hence we name this approach *iterative quadratic method (IQM)*. It is concluded as follows.

Algorithm 3.1. (IQM for Solving Problem (29))

1. Start: Given $A \in R^{m \times d}$, $Y \in R^{m \times n}$ and select parameters $\lambda > 0$, $q \in [1, 2]$ and $p \in (0, 2]$.
2. Set $k = 1$ and initialize $X_1 \in R^{d \times n}$.
3. For $k = 1, 2, \dots$ until convergence do :
 - $H_k = \text{diag}\{\frac{1}{\|X_k^i\|_2^{2-p}}\}_{i=1}^d$ ($0 < p < 2$) or $H_k = I_d$ ($p = 2$);
 - $C_k = -Y$;
 - For $i = 1 : I$
 - $B_i = A_i(X_k)_i$;
 - $C_k = B_i + C_k$;
 - end
 - $G_k = \text{diag}\{\frac{1}{\|C_k^i\|_2^{2-q}}\}_{i=1}^m$ ($1 \leq q < 2$) or $G_k = I_m$ ($q = 2$);
 - $M_k = A^T G_k A + \lambda \frac{p}{q} H_k$;
 - $X_{k+1} = M_k^{-1} A^T G_k Y$.

It is noticed that each iteration has to compute the inverse of M_k in Algorithm 3.1 which is expensive and unstable. Here we suggest to employ the general Penrose inverse of M_k to update the X_{k+1} . Moreover, the main computation $A_i X_i^*$ for classification is a by-product of B_i in computing the approximate solution X^* . Hence identifying test images can be achieved with minor extra calculations.

Algorithm 3.1 is a unified method solving $l_{2,q} - l_{2,p}$ -minimizations for $q \in [1, 2]$ and $p \in (0, 2]$. This approach provides algorithmic support to adaptively choose better fidelity measurement and regularization in various applications. Especially IQM provides a uniform algorithm for solving the existed representation based models: sparse representation ($q = 2$, $p = 1$), collaborative representation ($q = p = 2$) and l_1 -norm face recognition ($q = p = 1$).

3.2 Convergence Analysis of IQM

In this part, we will demonstrate the theoretical convergence of Algorithm 3.1. The key point is that the objective function $J(X)$ strictly decreases with respect to iterations until the matrix sequence $\{X_k\}$ converges to a stationary point of $J(X)$.

Lemma 3.1. Let $\varphi(t) = t - at^{\frac{1}{a}}$, where $a \in (0, 1)$. Then for any $t > 0$, $\varphi(t) \leq 1 - a$, and $t = 1$ is the unique maximizer.

Proof Taking the derivative of $\varphi(t)$ and set to zero, that is

$$\varphi'(t) = 1 - t^{\frac{1}{a}-1} = 0,$$

then $\varphi'(t) = 0$ has the unique solution $t = 1$ for any $a \in (0, 1)$ which is just the maximizer of $\varphi(t)$ in $(0, +\infty)$. \square

Lemma 3.2. Given X_k and X_{k+1} in $R^{d \times n}$, the following inequalities hold,

$$\|AX_{k+1} - Y\|_{2,q}^q - \frac{q}{2} \sum_{i=1}^m \frac{\|(AX_{k+1} - Y)^i\|_2^2}{\|(AX_k - Y)^i\|_2^{2-q}} \leq (1 - \frac{q}{2}) \|AX_k - Y\|_{2,q}^q \quad (33)$$

and

$$\|X_{k+1}\|_{2,p}^p - \frac{p}{2} \sum_{i=1}^d \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^{2-p}} \leq (1 - \frac{p}{2}) \|X_k\|_{2,p}^p \quad (34)$$

for any $q \in [1, 2)$ and $p \in (0, 2)$. Moreover, the equalities in Eq. (33) and (34) hold if and only if $\|(AX_{k+1} - Y)^i\|_2 = \|(AX_k - Y)^i\|_2$ for $i = 1, 2, \dots, m$ and $\|X_{k+1}^i\|_2 = \|X_k^i\|_2$ for $i = 1, 2, \dots, d$.

Proof Substituting $t_1 = \frac{\|(AX_{k+1} - Y)^i\|_2^q}{\|(AX_k - Y)^i\|_2^q}$ and setting $a_1 = \frac{q}{2}$ in Lemma 3.1, we obtain

$$\frac{\|(AX_{k+1} - Y)^i\|_2^q}{\|(AX_k - Y)^i\|_2^q} - \frac{q}{2} \frac{\|(AX_{k+1} - Y)^i\|_2^2}{\|(AX_k - Y)^i\|_2^2} \leq 1 - \frac{q}{2}. \quad (35)$$

Similarly taking $t_2 = \frac{\|X_{k+1}^i\|_2^p}{\|X_k^i\|_2^p}$ and $a_2 = \frac{p}{2}$ in $\varphi(t)$, we have

$$\frac{\|X_{k+1}^i\|_2^p}{\|X_k^i\|_2^p} - \frac{p}{2} \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^2} \leq 1 - \frac{p}{2}. \quad (36)$$

Multiplying Eq. (35) and Eq. (36) by $\|(AX_k - Y)^i\|_2^q$ and $\|X_k^i\|_2^p$ respectively, we have the following inequalities simultaneously

$$\|(AX_{k+1} - Y)^i\|_2^p - \frac{q}{2} \frac{\|(AX_{k+1} - Y)^i\|_2^2}{\|(AX_k - Y)^i\|_2^{2-q}} \leq (1 - \frac{q}{2}) \|(AX_k - Y)^i\|_2^q \quad (37)$$

for $i = 1, 2, \dots, m$, and

$$\|X_{k+1}^i\|_2^p - \frac{p}{2} \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^{2-p}} \leq (1 - \frac{p}{2}) \|X_k^i\|_2^p, \quad i = 1, 2, \dots, d. \quad (38)$$

Summing up i in formulas (37) and (38), we can derive (33) and (34).

Based on Lemma 3.1, $t_1 = 1$ and $t_2 = 1$ are the unique minimizers for $\varphi(t)$ in $(0, +\infty)$ when $a_1 = \frac{q}{2}$ and $a_2 = \frac{p}{2}$ respectively. Namely, $\|(AX_{k+1} - Y)^i\|_2 = \|(AX_k - Y)^i\|_2$ and $\|X_{k+1}^i\|_2 = \|X_k^i\|_2$ are necessary and sufficient for equalities hold in (37) and (38) respectively. \square

Remark 3.1. (33) and (34) are established nothing to do with Algorithm 3.1. The inequalities express the innate properties of mixed matrix norms $l_{2,q}$ - $l_{2,p}$ for $q \in [1, 2)$ and $p \in (0, 2)$.

Theorem 3.1. Suppose that $\{X_k\}$ is the matrix sequence generated by Algorithm 3.1. Then $J(X_k)$ strictly decreases with respect to k for any $1 \leq q \leq 2$ and $0 < p \leq 2$ until $\{X_k\}$ converges to a stationary point of $J(X)$.

Proof Based on the procedure of Algorithm 3.1, X_{k+1} is the solution to linear system (32), also the optimal matrix of problems (30) and (31). Thus we have

$$Q_k(X_{k+1}) \leq Q_k(X_k). \quad (39)$$

For $q \in [1, 2)$ and $p \in (0, 2)$, (39) is equivalent to

$$q \sum_{i=1}^m \frac{\|(AX_{k+1} - Y)^i\|_2^2}{\|(AX_k - Y)^i\|_2^{2-q}} + \lambda p \sum_{i=1}^d \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^{2-p}} \leq q \|AX_k - Y\|_{2,q}^q + \lambda p \|X_k\|_{2,p}^p, \quad (40)$$

It is noticed that $J(X_k) = \|AX_k - Y\|_{2,p}^p + \lambda \|X_k\|_{2,p}^p$. Adding inequalities (33) and $\lambda(34)$, the following formula will be derived

$$\begin{aligned} J(X_{k+1}) &- \left(\frac{q}{2} \sum_{i=1}^m \frac{\|(AX_{k+1} - Y)^i\|_2^2}{\|(AX_k - Y)^i\|_2^{2-q}} + \lambda \frac{p}{2} \sum_{i=1}^d \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^{2-p}} \right) \\ &\leq J(X_k) - \left(\frac{q}{2} \|AX_k - Y\|_{2,q}^q + \lambda \frac{p}{2} \|X_k\|_{2,p}^p \right). \end{aligned} \quad (41)$$

Based on (40) and (41), $J(X_{k+1}) \leq J(X_k)$ can be easily derived for $q \in [1, 2)$ and $p \in (0, 2)$.

For $q = 2$ or $p = 2$, the inequalities is much easier to derive. Taking $q = 2$ and $p \in (0, 2)$ for example, (39) is reduced to

$$\|AX_{k+1} - Y\|_{2,2}^2 + \lambda \frac{p}{2} \sum_{i=1}^d \frac{\|X_{k+1}^i\|_2^2}{\|X_k^i\|_2^{2-p}} \leq \|AX_k - Y\|_{2,2}^2 + \lambda \frac{p}{2} \|X_k\|_{2,p}^p, \quad (42)$$

Combining the formulas (42) and (34), we also obtain $J(X_{k+1}) \leq J(X_k)$. In the case of $q = 2$, $p \in (0, 2)$ or $q = p = 2$, $J(X_{k+1}) \leq J(X_k)$ can be deduced analogously.

Once $J(X_{k+1}) = J(X_k)$ happens for some k , the equalities in (40) and (41) (or (42)) hold. Hence the equalities in (33) and (34) are active. From Lemma 3.2, we obtain $\|(AX_{k+1} - Y)^i\|_2 = \|(AX_k - Y)^i\|_2$ for $i = 1, 2, \dots, m$ and $\|X_{k+1}^i\|_2 = \|X_k^i\|_2$ for $i = 1, 2, \dots, d$. Thus $G_{k+1} = G_k$ and $H_{k+1} = H_k$ which implies that X_{k+1} is a solution to (28). \square

The objective function sequence $\{J(X_k)\}$ is decreasing and lower bounded. Hence $\{J(X_k)\}$ eventually converges to some minimum of problem (16). The descending quantity measures the convergence precision.

Remark 3.2. The stopping criterion of Algorithm 3.1 can be chosen as $J(X_k) - J(X_{k+1}) \leq \epsilon$ or $\rho_k := \frac{J(X_k) - J(X_{k+1})}{J(X_k)} \leq \epsilon$ for some required precision $\epsilon > 0$.

Theoretically, $X_k^i = 0$ or $C_k^i = 0$ likely occurs in some step k , then H_k and G_k can not be well updated for non-Frobenius norm case ($0 < p < 2$ and $1 \leq q < 2$). We deal with it by perturbing with $\delta > 0$ such that $\{H_k\}_{ii} = \delta^{p-2} > 0$ and $\{G_k\}_{ii} = \delta^{q-2} > 0$. The descending of $\{J(X_k)\}$ is relaxed to

$$J(X_{k+1}) \leq J(X_k) + (1 - \frac{p}{2})\delta^p \quad \text{or} \quad J(X_{k+1}) \leq J(X_k) + (1 - \frac{q}{2})\delta^q. \quad (43)$$

If the convergence precision ϵ is chosen fairly larger than perturbation δ ($\epsilon \gg \delta$), perturbed $J(X_k)$ can be still considered approximate decreasing. As a matter of fact, $X_k^i = 0$ and $C_k^i = 0$ never happen in practical implementation.

4 Practical Implementation of JRC

In Algorithm 3.1, IQM has to update the matrix sequence by computing the inverse matrix of M_k . It is expensive in practical implementation especially for large scale problems. Reviewing the procedure of Algorithm 3.1, we notice that $X_{k+1} = M_k^{-1} A^T G_k Y$ exactly solves the k -th subproblem (31) which is unnecessary. It is observed that (31) is a quadratic positive definite subproblem. There are a lot of efficient algorithms to solve it approximately, such as conjugate gradient method, gradient methods with different stepsizes, etc. In this paper, we choose Barzilai and Borwein (BB) gradient method due to its simplicity and efficiency. BB gradient method was firstly presented in [18], afterwards extended and developed in many occasions and applications [18–23]. When applied to quadratic matrix optimization subproblem (31), the Barzilai and Borwein gradient method takes on

$$X_k^{(t+1)} = X_k^{(t)} - \alpha_k^{(t)} \nabla Q_k(X_k^{(t)}), \quad (44)$$

where the superscript (t) denotes the t -th iteration solving (31). $\nabla Q_k(X_k^{(t)})$ is the gradient matrix of $Q_k(X)$ with respect to $X_k^{(t)}$

$$\nabla Q_k(X_k^{(t)}) = M_k X_k^{(t)} - A^T G_k Y. \quad (45)$$

The Barzilai and Borwein gradient method [18] chose the stepsize $\alpha_k^{(t)}$ such that $D_k^{(t)} = \alpha_k^{(t)} I$ has a certain quasi-Newton property

$$D_k^{(t)} = \arg \min_{D=\alpha I} \|S_k^{(t-1)} - DT_k^{(t-1)}\|_F \quad (46)$$

or

$$D_k^{(t)} = \arg \min_{D=\alpha I} \|D^{-1} S_k^{(t-1)} - T_k^{(t-1)}\|_F, \quad (47)$$

where $\|\cdot\|_F$ denotes Frobenius matrix norm and $S_k^{(t-1)}$, $T_k^{(t-1)}$ are determined by the information achieved at the points $X_k^{(t)}$ and $X_k^{(t-1)}$

$$\begin{aligned} S_k^{(t-1)} &:= X_k^{(t)} - X_k^{(t-1)}; \\ T_k^{(t-1)} &:= \nabla Q_k(X_k^{(t)}) - \nabla Q_k(X_k^{(t-1)}) = M_k S_k^{(t-1)}. \end{aligned} \quad (48)$$

Solving (46) yields two BB stepsizes

$$\alpha_k^{(t)} = \frac{\text{Tr}((S_k^{(t-1)})^T T_k^{(t-1)})}{\text{Tr}((T_k^{(t-1)})^T T_k^{(t-1)})} \quad (49)$$

and

$$\alpha_k^{(t)} = \frac{\text{Tr}((S_k^{(t-1)})^T S_k^{(t-1)})}{\text{Tr}((S_k^{(t-1)})^T M_k S_k^{(t-1)})}. \quad (50)$$

Compared with the classical steepest descent method, BB gradient method often needs less computations but converges more rapidly [24]. For optimization problems

higher than two dimensions, BB method has theoretical difficulties due to its heavy non-monotone behavior. But for strongly convex quadratic problem with any dimension, BB method is convergent at R -linear rate [19, 21]. BB method has also been applied to matrix optimization problem [25] and exhibited desirable performance. Based on equations (44)-(50), the last step in Algorithm 3.1, $X_{k+1} = M_k^{-1} A^T G_k Y$, can be practically substituted by the BB gradient method as the k -th inner loop.

Algorithm 4.1. (*BB Gradient Method for Solving Subproblem (31)*)

1. *Start: given the inner loop stopping criterion $\epsilon_2 > 0$*
2. *Initialize $X_k^{(1)} = X_k$ and $\nabla Q_k^{(1)} = M_k X_k^{(1)} - A^T G_k Y$;*
3. *For $t = 1, 2, \dots$ until $\text{Tr}(\nabla Q_k^{(t)}) \leq \epsilon_2$, output $X_{k+1} = X_k^{(t)}$, do :*
 - if $t = 1$*
 - $$\alpha_k^{(t)} = \frac{\text{Tr}((\nabla Q_k^{(t)})^T \nabla Q_k^{(t)})}{\text{Tr}((\nabla Q_k^{(t)})^T M_k \nabla Q_k^{(t)})};$$
 - else*
 - $$S_k^{(t-1)} = X_k^{(t)} - X_k^{(t-1)};$$
 - $$T_k^{(t-1)} = \nabla Q_k^{(t)} - \nabla Q_k^{(t-1)};$$
 - $$\alpha_k^{(t)} \text{ is computed as (49) or (50);}$$
 - end*
 - $$X_k^{(t+1)} = X_k^{(t)} - \alpha_k^{(t)} \nabla Q_k^{(t)};$$
 - $$\nabla Q_k^{(t+1)} = M_k X_k^{(t+1)} - A^T G_k Y;$$

In the k -th inner loop, Algorithm 4.1 chooses two initial matrices. One is the approximate solution X_k to the last subproblem and another one is the Cauchy point from X_k [26]. The Cauchy stepsize $\alpha_k^{(1)}$ is the solution to the one-dimensional optimization problem

$$\min_{\alpha > 0} \phi(\alpha) := Q_k(X_k - \alpha \nabla Q_k(X_k)), \quad (51)$$

then the Cauchy point is $X_k + \alpha_k^{(1)} \nabla Q_k(X_k)$. If M_k in Algorithm 3.1 is guaranteed to be positive definite (if not, H_k or G_k can be slightly perturbed), subproblem (31) is a strongly convex quadratic. BB gradient method with step length (49) or (50) will converges at R -linear rate.

For simplicity, we name the IQM with inexact Algorithm 4.1 practically iterative quadratic method (PIQM). Still denote $\{X_k\}$ the approximate matrix sequence generated by PIQM. BB inner loop makes the objective function value of subproblem (31) decline, that is $Q(X_{k+1}) \leq Q(X_k)$. Then $\{J(X_k)\}$ is always decreasing which is sufficient and necessary for $\{X_k\}$ uniformly converging to the stationary point of problem (16). The following conclusion can be easily derived.

Theorem 4.1. *Denotes X^* the output point generated by PIQM, then X^* is an approximate stationary point of $J(X)$. Especially for $q, p \in [1, 2]$, X^* is an approximate global minimizer of optimization problem (16). When p is fractional, X^* is one of KKT points.*

An practical version of iteratively quadratic method for joint classification in face recognition can be concluded as follows.

Algorithm 4.2. (*PIQM for JRC*)

1. *Start: loading A, Y and setting $\lambda > 0, q \in [1, 2], p \in (0, 2]$ and precision levels $\varepsilon_1 > 0, \varepsilon_2 > 0$.*
2. *Employing PIQM to solve (16), output an approximate coding matrix $X^* := X_{k+1}$.*
3. *Classifying Y by X^* .*

5 Experimental Results

In this section, the joint representation based classification (JRC) with PIQM will be applied to face recognition. Three public data sets are used. Brief description is given as follows.

AT&T database is formerly known “the ORL database of faces”. It consists of 400 frontal images for 40 individuals. For each subject, 10 pictures were taken at different times, with varying lighting conditions, multiple facial expression, adornments and rotations up to 20 degree. All the images are aligned with dimension 112×92 . The database can be retrieved from http://www.cl.cam.ac.uk/Research/DTG/attarchive:pub/data/att_faces.tar.Z as a 4.5Mbyte compressed tar file. Typical pictures can be seen in Figure 1.



Figure 1: Typical images of AT & T database

Georgia-Tech database contains 15 images each of 50 subjects. The images are taken in two or three sessions at different times with different facial expressions, scale and background. The average size of the faces in these images is 150×150 pixels. Georgia Tech face database and the annotation can be found in http://www.ane.fian.com/research/face_r.eco.htm. Typical pictures of four persons are shown in Figure 2.



Figure 2: Typical images of Georgia-Tech database

Extended Yale B database consists of 2414 frontal-face images of 38 subjects. Each subject has around 64 images. The images are cropped and normalized to 192×168 under various laboratory-controlled lighting conditions [27, 28]. Figure 3 displays typical pictures of 4 subjects.

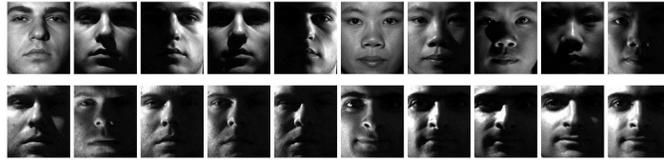


Figure 3: Typical images of Extended Yale B database

Extensive experiments are conducted for different image sizes and different parameters. Four comparable schemes are implemented, JRC, SRC, CRC-RLS and traditional SVM classifier. JRC is practically carried out via PIQM while SRC is solved by $l_1 - l_s$ solver [9] and CRC-RLS employs the code from [2]. We realize SVM by the software LIBSVM [30] with linear kernel, the pseudo code can be found in [http : //www.csie.ntu.edu.tw/ cjlin/libsvm/faq.html#f203](http://www.csie.ntu.edu.tw/~cjlin/libsvm/faq.html#f203). All the schemes are implemented by Matlab R2014a(win32) on a typical 4GiB memory and 2.40GHz PC.

Considering that JRC is a joint framework including SRC and CRC-RLS, we select six pairs of q, p in $[1, 2]$ and $(0, 2]$ respectively:

$$\begin{aligned} q = p = 2 & \text{ (corresponding to CRC-RLS),} \\ q = 2, p = 1 & \text{ (corresponding to SRC),} \end{aligned}$$

and other four generalized cases

$$\begin{aligned} q = 1.5 \ \& \ p = 1, \quad q = 1.5 \ \& \ p = 0.5, \\ q = 1 \ \& \ p = 1, \quad q = 1 \ \& \ p = 0.5. \end{aligned}$$

The parameter λ in (16) is varied from 0.01 to 10 each 10 times, and the best result is picked out. All the stopping precisions are set 10^{-3} .

All the images are re-sized like that of [1, 2]. For AT&T database, the pictures are down sampled to 11×10 . The downsampling ratios of Georgia-Tech database and Extended Yale B database are $1/8$ and $1/16$. For each subject, around 80% pictures are randomly selected for training and the left for testing. For example, 8 pictures of each individual in AT&T database are randomly picked out for training while the left 2 are for testing. All the classification schemes are directly applied to the images without any pre-processing. The recognition accuracy and running time are reported in Table 1-3.

Based on the experimental results on three databases, we draw the following conclusions:

- Jointly representing all the testing images simultaneously does accelerate face recognition. On all the databases, JRC ($q=p=2$) is the fastest one. The CPU time is thousand times less than that of SRC. For example, JRC ($q=p=2$) classifies 484 images in

Methods	The recognition accuracy	CPU time
SRC	98.75	67.2658
JRC(q=2,p=1)	97.5	0.1612
CRC-RLS	95	0.0872
JRC(q=p=2)	97.5	0.0073
SVM	95	0.0667
JRC(q=1.5,p=1)	97.5	0.3867
JRC(q=1.5,p=0.5)	95	1.8756
JRC(q=p=1)	97.5	0.1994
JRC(q=1,p=0.5)	97.5	0.1640

Table 1: The recognition accuracy (%) and running time (second) for AT&T database

Methods	Downsampling ratio 1/8		Downsampling ratio 1/16	
	Accuracy	Time	Accuracy	Time
SRC	99.33	2843	97.33	3197
JRC(q=2,p=1)	99.33	2.41	97.33	1.07
CRC-RLS	98	1.95	96.67	0.66
JRC(q=p=2)	99.33	0.97	98.67	0.17
SVM	96.67	5.09	96.67	1.46
JRC(q=1.5,p=1)	99.33	4.89	98.67	3.86
JRC(q=1.5,p=0.5)	99.33	4.89	98.67	3.89
JRC(q=p=1)	99.33	5.54	99.33	1.11
JRC(q=1,p=0.5)	99.33	4.79	99.33	1.09

Table 2: The recognition accuracy (%) and CPU time (second) for Georgia-Tech database

Methods	Down sampling ratio 1/8		Down sampling ratio 1/16	
	Accuracy	Time	Accuracy	Time
SRC	96.76	4828	96.36	668.53
JRC(q=2,p=1)	96.96	22.67	76.11	164.71
CRC-RLS	96.76	2.02	95.55	1.9
JRC(q=p=2)	96.96	0.75	91.29	0.34
SVM	95.55	6.12	94.33	2.61
JRC(q=1.5,p=1)	96.96	22.04	87.05	22.03
JRC(q=1.5,p=0.5)	96.96	54.21	65.59	101.59
JRC(q=p=1)	96.96	27.08	90.49	20.51
JRC(q=1,p=0.5)	96.96	26.87	91.29	25.23

Table 3: The recognition accuracy (%) and CPU time (second) for Extended Yale B database

0.17 second on Georgia-Tech database with downsampling ratio 1/16. And the accuracy rate is 98.67%, outperforming SRC (97.33%), CRC-RLS (96.67%) and SVM (96.67%). More details can be found in Table 1-3.

- JRC exhibits competitive performance in recognition accuracy. On AT & T database, the recognition rate of JRC is 97.5%, compared to 98.75% for SRC, 95% for CRC-RLS and SVM. On Georgia-Tech database, JRC achieves the best recognition rate (99.33%), consistently exceeds other classification schemes. On Yale B database with downsampling ratio 1/8, JRC also outperforms other methods in accuracy. Unfortunately, JRC does not keep the best achievement on downsampling ratio 1/16. The possible reason is that some pictures with strong contrast of lighting (see Figure 3) aggravates the noise for other images in joint coding.

- Different $q \in [1, 2]$ and $p \in (0, 2]$ for JRC indicate different feature pattern behind in the image set. Taking JRC ($q = 2, p = 1$) for example, the joint model combines sparsity of representation and correlation of multiple images. The representation coefficients reveal the joint effect on JRC ($q = 2, p = 1$), Figure 4 gives an example from Yale B database. Compared to SRC, JRC ($q = 2, p = 1$) concentrates a group sparsity but not a single one. Actually, the other testing samples (12 pictures) of the same subject also have the similar group representation pattern.

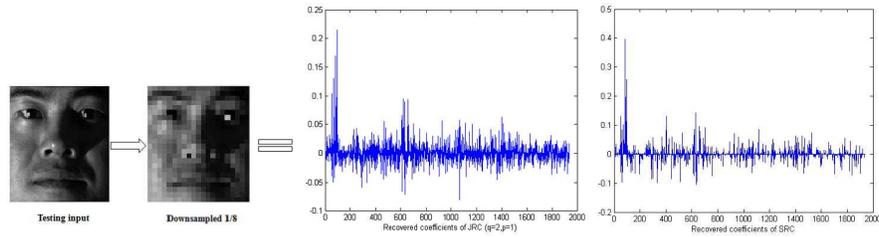


Figure 4: The recovered coefficients by JRC ($q=2,p=1$) and SRC

- The convergence behavior of PIQM for JRC is displayed in Figures 5. The x axis is the iterations and y-axis stands for the logarithm of ρ_k . PIQM converges within 40 steps on three databases for all jointly sparse models (five pairs q and p). JRC ($q=p=2$) always converges in three iterations hence its plot is omitted here. Anyway, PIQM provides a uniform algorithm for varied JRC with respect to $q \in [1, 2]$ and $p \in (0, 2]$.

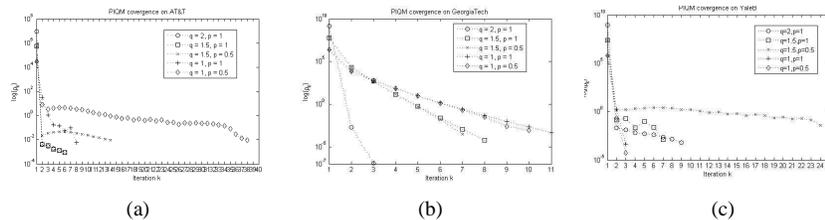


Figure 5: (a) PIQM on AT & T (b) PIQM on Georgia-Tech (c) PIQM on Yale B

• From Table 1-3, it is observed that CRC-RLS has a fairly good performance in recognition accuracy and CPU time. But CRC-RLS is heavily sensitive to the regularization parameter λ (see Table 4) because it has a smooth regularizer. By comparison, JRC ($q=p=2$) is more stable for its joint technique. Multiple images has complementary effect for recognition especially when the model is ill-posed.

$\lambda =$	0.01	0.1	1	10	100
CRC-RLS	28.34	66.82	95	96.76	96.76
JRC($q=p=2$)	96.96	96.96	96.96	96.96	96.96

Table 4: The recognition accuracy (%) for different λ on Extended Yale B database with downsampling ratio 1/8

6 Conclusions

In this paper, a joint representation classification for collective face recognition is proposed. By aligning all the testing images into a matrix, joint representation coding is reduced to a kind of generalized matrix pseudo norm based optimization problems. A unified algorithm is developed to solve the mixed $l_{2,q}-l_{2,p}$ -minimizations for $q \in [1, 2]$ and $p \in (0, 2]$. The convergence is also uniformly demonstrated. To adapt the algorithm to the large scale case, a practical iterative quadratic method is considered to inexactly solve the subproblems. Experiment results on three data-sets validate the collective performance of the proposed scheme. The joint representation based classification is confirmed to improve the performance in recognition rate and running time than the state-of-the-arts.

Acknowledgement The first author thanks software engineer Luo Aiwen for his code support.

References

- [1] J. Wright, A.Y. Yang, A. Ganesh, S.S. Sastry and Y. Ma. Robust face recognition via sparse representation. IEEE Trans. PAMI, 2009, (31)(2):210-227.
- [2] L. Zhang, M. Yang and X.C. Feng. Sparse representation or collaborative representation: which help face recognition? 13th International Conference on Computer Vision, 2011, pp. 471-478.
- [3] J. Wright and Y. Ma. Dense error correction via l_1 minimization. IEEE Transactions on Information Theory, 2010, 56(7):3540-3560.
- [4] S.H. Gao, I.W.H. Tsang and L.T. Chia. Kernel sparse representation for image classification and face recognition. In ECCV, 2010.

- [5] J.Z. Huang, X.L. Huang and D. Metaxas. Simultaneous image transformation and sparse representation recovery. In CVPR, 2008.
- [6] A. Wagner, J. Wright, W. Xu and Y. Ma. Towards a practical face recognition system: robust registration and illumination by sparse representation. In CVPR, 2009.
- [7] Y. Peng, A. Wagner, J. Wright, W. Xu and Y. Ma. RASL: robust alignment by sparse and low-rank decomposition for linearly correlated image. IEEE Trans. PAMI, 2012, 34(11):2233-2246.
- [8] R.P. Wang, S.G. Shan, X.L. Chen and W. Gao. Manifold-manifold method distance with application to face recognition based on image set. IEEE Transactions on Image Processing, 2012, 21(10):4466-4479.
- [9] S.J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. A interior-point method for large-scale ℓ_1 -regularized least squares. IEEE Journal on Selected Topics in Signal Processing, 2007, 1(4):606C617.
- [10] A. Rakotomamonjy, R. Flamary, G. Gasso, and S. Canu. $l_p - l_q$ Penalty for sparse linear and sparse multiple kernel multitask learning. IEEE Transactions on Neural Networks, 2011, (22)(8):1307-1320.
- [11] S. Suvrit. Fast projection onto $l_{1,q}$ -norm balls for grouped feature selection. In Proceeding of Machine Learning and Knowledge Discovery in Databases, 2011, Athens, Greece.
- [12] S. Sumit, M.P. Vishal, M.N. Nasser, and C. Rama. Joint sparse representation for robust multimodal bimetrics recognition. IEEE Trans. PAMI, 2014, (36)(1):113-126.
- [13] R. Chartrand and W.T. Yin. Iteratively reweighed algorithms for compressive sensing. The 33rd International Conference on Acoustics, Speech, and Signal Processing, 2008, pp. 3869-3872.
- [14] R. Chartrand. Exact reconstructions of sparse signals via nonconvex minimization. IEEE Signal Processing Letters, 2007, 14(10):707-710.
- [15] Z.B. Xu, H. Zhang, Y. Wang, X.Y. Chang and Y. Liang. $L_{\frac{1}{2}}$ regularizer. Science in China: Series F, 2010, 52(6):1159-1169.
- [16] H. Wang, F.P. Nie, W.D. Cai and H. Huang. Semi-supervised robust dictionary learning via efficient $l_{2,0+}$ -norms minimizations. IEEE International Conference on Computer Vision, 2013, pp. 1145-1152.
- [17] L.P. Wang, S.C. Chen and Y.P. Wang. A unified algorithm for mixed $l_{2,p}$ -minimizations and its application in feature selection. Computational Optimization and Applications, 2014, 58:409-421.
- [18] J. Barzilai and J.M. Borwein. Two-point step size gradient methods. IMA Journal of Numerical Analysis, 1988, 8:141-148.

- [19] M. Raydan. On the Barzilai and Borwein choice of steplength for the gradient method. *IMA Journal of Numerical Analysis*, 1993, 13:321-326.
- [20] M. Raydan. The Barzilai and Borwein gradient method for the large scale unconstrained minimization problem. *SIAM Journal on Optimization*, 1997, 7:26-33.
- [21] Y.H. Dai and L.Z. Liao. R -linear convergence of the Barzilai and Borwein gradient method. *IMA Journal of Numerical Analysis*, 2002, 26:1-10.
- [22] Y.H. Dai and R. Fletcher. New algorithms for singly linearly constrained quadratic programs subject to lower and upper bounds. *Math. Program. Ser. A*, 2006, 106:403-421.
- [23] Y.X. Yuan. A new stepsize for the steepest descent method. *Journal of Computational Mathematics*. 2006, 24(2):149-156.
- [24] R. Fletcher. Low storage method for unconstrained optimization. *Lectures Appl. Math. (AMS)*, 1990, 26:165-179.
- [25] B. Jiang, C.F. Cui and Y.H. Dai. Unconstrained optimization models for computing several extreme eigenpairs of real symmetric matrices. *Pacific Journal of Optimization*. 2014, 10(1):55 - 71.
- [26] A. Cauchy. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comp. Rend. Sci. Pari*, 1847, 25:141-148.
- [27] A. Georghiades, P. Belhumeur and D. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. PAMI*, 2001, 23(6):643-660.
- [28] L. Lee, J. Ho, and D. Kriegman. Acquiring linear subspaces for face recognition under variable lighting. *IEEE Trans. PAMI*, 2005, 27(5):684-698.
- [29] A. Martinez and R. Benavente. The AR face database. *CVC Tech. Report No. 24*, 1998.
- [30] R.E. Fan, P.H. Chen and C.J. Lin. Working set selection using second order information for training SVM. *Journal of Machine Learning Research* 2005, 6:1889-1918.