

# Sparse Representation or Collaborative Representation: Which Helps Face Recognition?

Lei Zhang<sup>a</sup>, Meng Yang<sup>a</sup>, and Xiangchu Feng<sup>b</sup>

<sup>a</sup>Dept. of Computing, The Hong Kong Polytechnic University, Hong Kong, China

<sup>b</sup>Dept. of Applied Mathematics, Xidian University, Xi'an, China  
{cslzhang, csmyang}@comp.polyu.edu.hk

## Abstract

*As a recently proposed technique, sparse representation based classification (SRC) has been widely used for face recognition (FR). SRC first codes a testing sample as a sparse linear combination of all the training samples, and then classifies the testing sample by evaluating which class leads to the minimum representation error. While the importance of sparsity is much emphasized in SRC and many related works, the use of collaborative representation (CR) in SRC is ignored by most literature. However, is it really the  $l_1$ -norm sparsity that improves the FR accuracy? This paper devotes to analyze the working mechanism of SRC, and indicates that it is the CR but not the  $l_1$ -norm sparsity that makes SRC powerful for face classification. Consequently, we propose a very simple yet much more efficient face classification scheme, namely CR based classification with regularized least square (CRC\_RLS). The extensive experiments clearly show that CRC\_RLS has very competitive classification results, while it has significantly less complexity than SRC.*

## 1. Introduction

It has been found that natural images can be sparsely coded by structural primitives [1], and in recent years sparse coding or sparse representation has been widely studied to solve the inverse problems in various image restoration applications [2-3], partially due to the progress of  $l_0$ -norm and  $l_1$ -norm minimization techniques [4-6].

Recently, sparse representation has also been used in pattern classification. Huang *et al.* [7] sparsely coded a signal over a set of redundant bases and classified the signal based on its coding vector. In [8], Wright *et al.* reported a very interesting work by using sparse representation for robust face recognition (FR). A query face image is first sparsely coded over the template images, and then the classification is performed by checking which class yields the least coding error. Such a sparse representation based classification (SRC) scheme achieves a great success in FR, and it boosts the research of sparsity based pattern classification. Gao *et al.* [9] proposed the kernel sparse

representation for FR, while Yang and Zhang [10] used the Gabor features for SRC with a learned Gabor occlusion dictionary to reduce the computational cost. Cheng *et al.* [11] discussed the  $l_1$ -graph for classification, and Yang *et al.* [12] combined sparse coding with linear spatial pyramid matching for image classification. A recent review of sparse representation for computer vision and pattern recognition applications can be found in [13].

In sparse representation based FR, usually we assume that the face images are aligned. Recently, sparse representation has been extended to solve the misalignment or pose change. The method in [14] is invariant to image-plane transformation. The method in [15] could deal with misalignment and illumination variation. In [16], Peng *et al.* studied how to simultaneously align a batch of linearly correlated images with gross corruption.

Sparse representation (or coding) codes a signal  $\mathbf{y}$  over a dictionary  $\Phi$  such that  $\mathbf{y} \approx \Phi \boldsymbol{\alpha}$  and  $\boldsymbol{\alpha}$  is a sparse vector. The sparsity of  $\boldsymbol{\alpha}$  can be measured by  $l_0$ -norm, which counts the number of non-zeros in  $\boldsymbol{\alpha}$ . Since the combinatorial  $l_0$ -minimization is NP-hard, the  $l_1$ -minimization, as the closest convex function to  $l_0$ -minimization, is widely employed in sparse coding:  $\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_1$  s.t.  $\|\mathbf{y} - \Phi \boldsymbol{\alpha}\|_2 \leq \varepsilon$ , where  $\varepsilon$  is a small constant. Although  $l_1$ -minimization is much more efficient than  $l_0$ -minimization, it is still time consuming, and hence many fast algorithms were proposed to speed up the  $l_1$ -minimization process.

As reviewed in [17], there are five representative fast  $l_1$ -minimization approaches: Gradient Projection, Homotopy, Iterative Shrinkage-Thresholding, Proximal Gradient, and Augmented Lagrange Multiplier (ALM). It was indicated that for noisy data, first order  $l_1$ -minimization techniques (e.g., SpaRSA [18], FISTA [19] and ALM [20]) are more efficient, while for FR, Homotopy [21], ALM and  $l_1$ -ls [22] are better for their accuracy and fast speed.

Although SRC [8] has shown interesting results in FR and has been widely studied in the community, its working mechanism has not been clearly revealed yet. Most literature, including [8], emphasizes too much on the role of  $l_1$ -norm sparsity in face classification, while the role of collaborative representation (CR), i.e., using the training samples from all classes to represent the query sample  $\mathbf{y}$ , is much ignored. The  $l_1$ -minimization makes the sparsity

based classification schemes such as SRC very expensive; however, is it really the  $l_1$ -norm sparsity that makes SRC powerful for FR? Very recently some researchers have started to question the use of sparsity in image classification, such as [29-30].

This paper devotes to analyze the working mechanism of SRC. We will explain why sparsity could improve discrimination, and more importantly, we will indicate that it is the CR, but not the  $l_1$ -norm sparsity, that plays the essential role for classification in SRC. Consequently, we propose a new classification scheme, namely CR based classification with regularized least square (CRC\_RLS), which has significantly less complexity than SRC but leads to very competitive classification results.

Section 2 briefly reviews SRC. Section 3 analyzes sparse representation and CR. Section 4 presents the CRC\_RLS scheme. Section 5 conducts extensive experiments, and Section 6 concludes the paper.

## 2. The SRC scheme

**Table 1:** The SRC Algorithm

1. Normalize the columns of  $X$  to have unit  $l_2$ -norm.

2. Code  $y$  over  $X$  via  $l_1$ -minimization

$$(\hat{\alpha}) = \arg \min_{\alpha} \|\alpha\|_1 \text{ s.t. } \|y - X\alpha\|_2 < \varepsilon \quad (1)$$

where constant  $\varepsilon$  is to account for the dense small noise in  $y$ , or to balance the coding error of  $y$  and the sparsity of  $\alpha$ .

3. Compute the residuals

$$e_i(y) = \|y - X_i \hat{\alpha}_i\|_2 \quad (2)$$

where  $\hat{\alpha}_i$  is the coding coefficient vector associated with class  $i$ .

4. Output the identity of  $y$  as

$$\text{identity}(y) = \arg \min_i \{e_i\} \quad (3)$$

Denote by  $X_i \in \mathfrak{R}^{m \times n}$  the dataset of the  $i^{\text{th}}$  class, and each column of  $X_i$  is a sample of class  $i$ . Suppose that we have  $K$  classes of subjects, and let  $X = [X_1, X_2, \dots, X_K]$ . Once a query image  $y \in \mathfrak{R}^m$  comes, we code it as  $y \approx X\alpha$ , where  $\alpha = [\alpha_1; \dots; \alpha_i; \dots; \alpha_K]$  and  $\alpha_i$  is the coding vector associated with class  $i$ . If  $y$  is from the  $i^{\text{th}}$  class, usually  $y \approx X_i \alpha_i$  holds well, implying that most coefficients in  $\alpha_k$ ,  $k \neq i$ , are nearly zeros and only  $\alpha_i$  has significant entries. That is, the sparse non-zero entries in  $\alpha$  can encode the identity of sample  $y$ . The procedures of SRC are summarized in Table 1.

## 3. Sparse representation and collaborative representation

From Table 1, we see that there are two key points in SRC. The first key point is that the coding vector of query sample  $y$  is required to be sparse, and the second key point

is that the coding of  $y$  is performed collaboratively over the whole dataset  $X$  instead of each subset  $X_i$ . Suppose that  $y$  belongs to some class in the dataset, it was claimed in [8] that the sparsest (or the most compact) representation of  $y$  over  $X$  is naturally discriminative and thus can indicate the identity of  $y$ . It was also claimed that SRC is a generalization of the classical nearest neighbor (NN) and nearest subspace (NS) classifiers. The NN classifier represents  $y$  by each individual of the training samples; the NS classifier represents  $y$  by the training samples of each class; and SRC represents  $y$  collaboratively by samples of all classes. In this section, we first illustrate why sparsity makes representation more discriminative, and then discuss the collaborative representation involved in SRC.

### 3.1. Why sparse representation?

Denote by  $\Phi \in \mathfrak{R}^{m \times n}$  a dictionary of atoms. If  $\Phi$  is complete, then any signal  $x \in \mathfrak{R}^m$  can be accurately represented as the linear combination of the atoms in  $\Phi$ . If  $\Phi$  is orthogonal, however, often we need to use many atoms from  $\Phi$  to faithfully represent  $x$ . If we want to use less atoms to represent  $x$ , we must relax the orthogonality imposed on  $\Phi$ . In other words, we must allow more atoms to be involved in  $\Phi$  so that we have more choices to represent  $x$ , leading to an over-complete dictionary  $\Phi$  but a sparser representation of signal  $x$ . For example, it is well-known that redundant wavelet transforms have much better denoising performance than orthogonal wavelet transforms. The great success of sparse representation in image restoration [2-3] further validates this.

In the scenario of FR, each class of face images often lies in a small subspace of  $\mathfrak{R}^m$ . That is, the  $m$ -dimensional face image  $x$  can be characterized by a feature vector of much lower dimensionality. If we take the set of training samples of class  $i$ , i.e.,  $X_i$ , as the dictionary for this class, in practice the atoms (i.e., the training samples) of  $X_i$  will be correlated. Assume that we have enough training samples for each class so that all the images of class  $i$  can be faithfully represented by  $X_i$ , then  $X_i$  is an over-complete dictionary<sup>1</sup> because of the correlation of training samples of class  $i$ , and we can conclude that a testing sample  $y$  of class  $i$  can be sparsely represented over dictionary  $X_i$ .

Another important fact in FR is that all the face images are somewhat similar, while some subjects may have very similar face images. This implies that dictionary  $X_i$  and dictionary  $X_j$  are not incoherent but can be highly correlated. Let  $X_j = X_i + \Delta$ . Using the NS classifier, for a query sample  $y$  from class  $i$ , we can calculate by least square method a vector  $\alpha_i = \arg \min_{\alpha} \|y - X_i \alpha\|_2$ . Let  $e_i =$

<sup>1</sup> More strictly speaking, it should be the dimensionality reduced dictionary of  $X_i$  that is over-complete. For the convenience of expression, we simply use  $X_i$  in the development.

$\mathbf{y}-\mathbf{X}_i\boldsymbol{\alpha}_i$ . Similarly, if we represent  $\mathbf{y}$  by class  $j$ , there is  $\boldsymbol{\alpha}_j = \arg \min_{\boldsymbol{\alpha}} \|\mathbf{y} - \mathbf{X}_j\boldsymbol{\alpha}\|_2$  and we let  $\mathbf{e}_j = \mathbf{y} - \mathbf{X}_j\boldsymbol{\alpha}_j$ . Suppose that  $\mathbf{X}_i, \mathbf{X}_j \in \mathbb{R}^{m \times n}$ , if  $\Delta$  is small such that

$$\xi = \frac{\|\Delta\|_F}{\|\mathbf{X}_i\|_F} \leq \frac{\sigma_n(\mathbf{X}_i)}{\sigma_1(\mathbf{X}_i)}$$

where  $\sigma_1(\mathbf{X}_i)$  and  $\sigma_n(\mathbf{X}_i)$  are the largest and smallest eigenvalues of  $\mathbf{X}_i$ , respectively, then we have the following relationship between  $\mathbf{e}_i$  and  $\mathbf{e}_j$  (page 242, [28]):

$$\frac{\|\mathbf{e}_j - \mathbf{e}_i\|_2}{\|\mathbf{y}\|_2} \leq \xi (1 + \kappa_2(\mathbf{X}_i)) \min\{1, m - n\} + O(\xi^2) \quad (4)$$

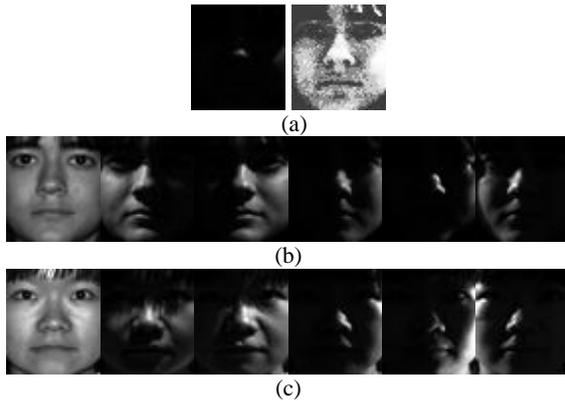
where  $\kappa_2(\mathbf{X}_i)$  is the  $l_2$ -norm conditional number of  $\mathbf{X}_i$ .

From Eq. (4), we can clearly see that if  $\Delta$  is small, i.e., subjects  $i$  and  $j$  look similar to each other, then the distance between  $\mathbf{e}_i$  and  $\mathbf{e}_j$  can be very small. This makes the classification very unstable because a small disturbance can lead to  $\|\mathbf{e}_j\|_2 < \|\mathbf{e}_i\|_2$ , resulting in a wrong classification.

The above problem can be much alleviated by imposing some sparsity on  $\boldsymbol{\alpha}_i$  and  $\boldsymbol{\alpha}_j$ . The reason is very simple. If  $\mathbf{y}$  is from class  $i$ , it is more likely that we can use only a few samples (e.g., 5 or 6 samples) in  $\mathbf{X}_i$  to represent  $\mathbf{y}$  with a good accuracy. In contrast, we may need more samples (e.g., 8 or 9 samples) in  $\mathbf{X}_j$  to represent  $\mathbf{y}$  with nearly the same accuracy. Under a certain sparsity constraint, the representation error of  $\mathbf{y}$  by  $\mathbf{X}_i$  will be visibly lower than that by  $\mathbf{X}_j$ , making the classification of  $\mathbf{y}$  easier. The sparse representation of  $\mathbf{y}$  by  $\Phi$  can be formulated as

$$\min_{\boldsymbol{\alpha}} \|\mathbf{y} - \Phi\boldsymbol{\alpha}\|_2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_p \leq \varepsilon \quad (5)$$

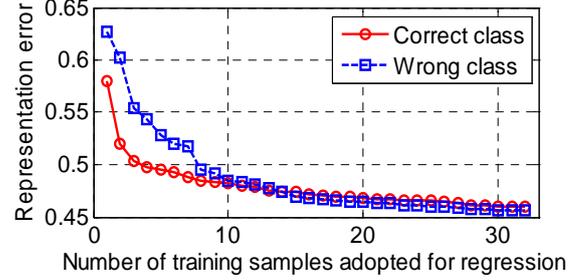
where  $\varepsilon$  is a constant and  $p$  can be 0, 1, or any other eligible sparsity metric.



**Figure 1:** (a) The query face image (left: original image, right: the one after histogram equalization for better visualization); (b) some training samples from the class of the query image; (c) some training samples from another class.

Let  $e = \|\mathbf{y} - \Phi\boldsymbol{\alpha}\|_2$  and set  $p=0$  in Eq. (5). We could plot the curve of  $e$  versus  $\varepsilon$  to illustrate why sparsity improves

discrimination. Fig. 1(a) shows a testing face image from class 32 in the Extended Yale B database. Some training samples of class 32 are shown in Fig. 1(b). Some training samples of class 5, which looks similar to class 32, are shown in Fig. 1(c). We use the training samples of the two classes as dictionaries to represent the query sample in Fig. 1(a), respectively, under different sparsity  $\varepsilon$ . The two “ $e$  vs.  $\varepsilon$ ” curves are drawn in Fig. 2.



**Figure 2:** The curve of representation error versus the number of training samples in each class.

From Fig. 2, we can see that when using only a few training samples ( $<3$ ) to represent the query sample, both of the two classes have big error; when more and more training samples are involved, the representation error decreases. However, the discrimination ability of the representation error will also reduce if too many samples ( $>10$ ) are used. Thus, the sparsity of coefficients should be considered. From the above analyses, we may have the following proposition for  $l_0$ -norm sparsity: *a query sample should be classified to the class which could faithfully represent it using less number of samples.*

In practice, if the number of training samples of each class is relatively big, we can represent the testing sample  $\mathbf{y}$  class by class. Since  $l_0$ -minimization is combinatorial NP-hard,  $l_1$ -minimization with the following Lagrangian formulation is often adopted:

$$(\boldsymbol{\alpha}_i) = \arg \min_{\boldsymbol{\alpha}} \left\{ \|\mathbf{y} - \mathbf{X}_i\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right\} \quad (6)$$

Both the representation error  $e_i = \|\mathbf{y} - \mathbf{X}_i\boldsymbol{\alpha}_i\|_2$  and the sparsity term  $\|\boldsymbol{\alpha}_i\|_1$  can be used to classify the sample. Later we will see that (please refer to the experimental results in Section 5.2) actually the  $l_2$ -norm can also be used to regularize  $\boldsymbol{\alpha}_i$ , and the  $l_1$ -norm and  $l_2$ -norm lead to almost the same result.

### 3.2. Why collaborative representation?

In our discussion in Section 3.1, we assumed that there are enough training samples for each class so that the (dimensionality reduced) dictionary  $\mathbf{X}_i$  is over-complete. Unfortunately, FR is a typical small-sample-size problem, and  $\mathbf{X}_i$  is under-complete in general. If we use  $\mathbf{X}_i$  to represent  $\mathbf{y}$ , the representation error can be big, even when

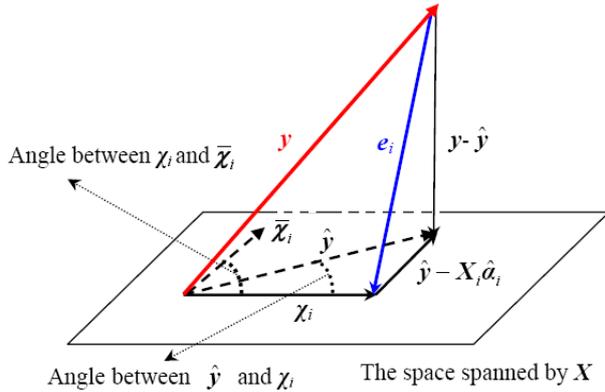
$\mathbf{y}$  is from class  $i$ . Consequently, the classification will be unstable not matter the error  $e_i$  or the sparsity  $\|\alpha_i\|_p$  or both of them are used for decision making.

One obvious solution to solving this problem is to use more samples of class  $i$  to represent  $\mathbf{y}$ . But where could we have these samples? Fortunately, one fact in FR is that face images of different classes share similarities. Some sample from class  $j$  may be very helpful to represent the testing sample with label  $i$ . In SRC [8], this ‘‘lack of samples’’ problem is solved by taking the face images from all the other classes as the possible samples of each class. That is, it codes the testing image  $\mathbf{y}$  collaboratively over the dictionary of all samples  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$  under the  $l_1$ -norm sparsity constraint.

One interesting point here is that after the collaborative representation (CR) with all classes, SRC classifies  $\mathbf{y}$  individually (i.e., check class by class). For the simplicity of analysis, let’s remove the  $l_1$ -norm sparsity term in Eq. (1), and then the representation becomes a least square problem:  $(\hat{\alpha}) = \arg \min_{\alpha} \|\mathbf{y} - \mathbf{X}\alpha\|_2^2$ . The associated representation  $\hat{\mathbf{y}} = \sum_i \mathbf{X}_i \hat{\alpha}_i$  is actually the perpendicular projection of  $\mathbf{y}$  onto the space spanned by  $\mathbf{X}$ . In SRC, the reconstruction error by each class  $e_i = \|\mathbf{y} - \mathbf{X}_i \hat{\alpha}_i\|_2^2$  is used for classification. It can be readily derived that

$$e_i = \|\mathbf{y} - \mathbf{X}_i \hat{\alpha}_i\|_2^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{y}} - \mathbf{X}_i \hat{\alpha}_i\|_2^2$$

Obviously, it is the amount  $e_i^* = \|\hat{\mathbf{y}} - \mathbf{X}_i \hat{\alpha}_i\|_2^2$  that works for classification because  $\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$  is a constant for all classes.



**Figure 3:** Geometric illustration of the representation of  $\mathbf{y}$  over  $\mathbf{X}$ .

Denote by  $\chi_i = \mathbf{X}_i \hat{\alpha}_i$  and  $\bar{\chi}_i = \sum_{j \neq i} \mathbf{X}_j \hat{\alpha}_j$ . Fig. 3 shows geometrically the representation of  $\mathbf{y}$  over  $\mathbf{X}$ . Since  $\bar{\chi}_i$  is parallel to  $\hat{\mathbf{y}} - \mathbf{X}_i \hat{\alpha}_i$ , we can readily have

$$\frac{\|\hat{\mathbf{y}}\|_2}{\sin(\chi_i, \bar{\chi}_i)} = \frac{\|\hat{\mathbf{y}} - \mathbf{X}_i \hat{\alpha}_i\|_2}{\sin(\hat{\mathbf{y}}, \chi_i)}$$

where  $(\chi_i, \bar{\chi}_i)$  is the angle between  $\chi_i$  and  $\bar{\chi}_i$ , and  $(\hat{\mathbf{y}}, \chi_i)$  is the angle between  $\hat{\mathbf{y}}$  and  $\chi_i$ . Finally, the representation error can be represented by

$$e_i^* = \frac{\sin^2(\hat{\mathbf{y}}, \chi_i) \|\hat{\mathbf{y}}\|_2^2}{\sin^2(\chi_i, \bar{\chi}_i)} \quad (7)$$

Eq. (7) shows that by using CR, when we judge if  $\mathbf{y}$  belongs to class  $i$ , we will not only consider if the angle between  $\hat{\mathbf{y}}$  and  $\chi_i$  is small (i.e., if  $\sin(\hat{\mathbf{y}}, \chi_i)$  is small), we will also consider if the angle between  $\chi_i$  and  $\bar{\chi}_i$  is big (i.e., if  $\sin(\chi_i, \bar{\chi}_i)$  is big). Such a ‘‘double checking’’ makes the classification more effective and robust.

One problem is that when the number of classes is too big, the least square solution  $(\hat{\alpha}) = \min_{\alpha} \|\mathbf{y} - \mathbf{X}\alpha\|_2^2$  will become unstable. In SRC, the  $l_1$ -norm sparsity constraint is imposed on  $\alpha$  to make the solution stable. However, it is not necessary to use the strong  $l_1$ -norm to this end. As we will see next, by using the much weaker  $l_2$ -norm to regularize the solution of  $\alpha$ , we can have similar classification results but with significantly lower complexity. *In summary, it is the CR but not the  $l_1$ -norm sparsity constraint that truly improves the FR performance.*

#### 4. Collaborative representation based classification (CRC)

Most of the previous works [8-13] emphasize the importance of sparsity for classification but do not investigate much the role of collaboration between classes in representing the query sample. Is it really the sparsity that improves the FR accuracy? Or is it the CR that truly helps FR? To answer this question, we propose here a simple CR based classification (CRC) scheme, and conduct experiments to give the answer in next section.

In order to collaboratively represent the query sample using  $\mathbf{X}$  with low computational burden, we propose to use the regularized least square method. There is

$$(\hat{\rho}) = \arg \min_{\rho} \left\{ \|\mathbf{y} - \mathbf{X} \cdot \rho\|_2^2 + \lambda \|\rho\|_2^2 \right\} \quad (8)$$

where  $\lambda$  is the regularization parameter. The role of the regularization term is twofold. First, it makes the least square solution stable, and second, it introduces a certain amount of ‘‘sparsity’’ to the solution  $\hat{\rho}$ , yet this sparsity is much weaker than that by  $l_1$ -norm.

The solution of CR with regularized least square in Eq. (8) can be easily and analytically derived as

$$\hat{\rho} = (\mathbf{X}^T \mathbf{X} + \lambda \cdot \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (9)$$

Let  $\mathbf{P} = (\mathbf{X}^T \mathbf{X} + \lambda \cdot \mathbf{I})^{-1} \mathbf{X}^T$ . Clearly,  $\mathbf{P}$  is independent of  $\mathbf{y}$  so that it can be pre-calculated as a projection matrix. Once a query sample  $\mathbf{y}$  comes, we can just simply project  $\mathbf{y}$  onto  $\mathbf{P}$

via  $\mathbf{P}\mathbf{y}$ . This makes CR very fast.

The classification by  $\hat{\boldsymbol{\rho}}$  is similar to the classification by  $\hat{\boldsymbol{\alpha}}$  in SRC (refer to Table 1). In addition to the class specific representation residual  $\|\mathbf{y} - \mathbf{X}_i \cdot \hat{\boldsymbol{\rho}}_i\|_2$ , where  $\hat{\boldsymbol{\rho}}_i$  is the coefficient vector associated with class  $i$ , the  $l_2$ -norm ‘‘sparsity’’  $\|\hat{\boldsymbol{\rho}}_i\|_2$  can also bring some discrimination information for classification. Therefore we propose to use both of them in classification. Based on our experiments, this improves slightly the classification accuracy over that by using only the representation residual. The proposed CRC with regularized least square (CRC\_RLS) algorithm is summarized as follows.

**Table 2:** The CRC\_RLS Algorithm

1. Normalize the columns of  $\mathbf{X}$  to have unit  $l_2$ -norm.
2. Code  $\mathbf{y}$  over  $\mathbf{X}$  by

$$\hat{\boldsymbol{\rho}} = \mathbf{P}\mathbf{y}$$

$$\text{where } \mathbf{P} = (\mathbf{X}^T \mathbf{X} + \lambda \cdot \mathbf{I})^{-1} \mathbf{X}^T.$$

3. Compute the regularized residuals

$$r_i = \|\mathbf{y} - \mathbf{X}_i \cdot \hat{\boldsymbol{\rho}}_i\|_2 / \|\hat{\boldsymbol{\rho}}_i\|_2 \quad (10)$$

4. Output the identity of  $\mathbf{y}$  as

$$\text{Identity}(\mathbf{y}) = \text{argmin}_i \{r_i\}.$$

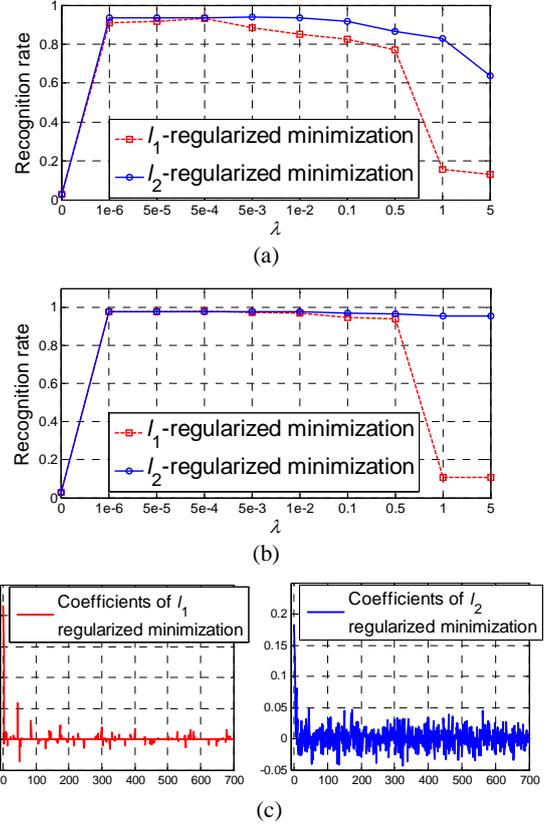
## 5. Experimental results

Considering the accuracy and efficiency, we chose  $l_1$ - $l_2$  [22] to solve the  $l_1$ -regularized minimization in SRC. All the experiments were carried out using MATLAB on a 3.16 GHz machine with 3.25GB RAM. In the experiments of gender classification in Section 5.2, the parameter  $\lambda$  in CRC\_RLS is set as 0.08. In FR, considering that when more classes (and thus more samples) are used, the least square solution of CR will be more unstable and thus higher regularization is required, we set  $\lambda$  as  $0.001 * n / 700$  in all FR experiments, where  $n$  is the number of training samples. The MATLAB code of CRC\_RLS can be downloaded at <http://www4.comp.polyu.edu.hk/~cslzhang/code.htm>.

Three face databases, including Extended Yale B [23] [24], AR [25], and large-scale Multi-PIE [26], are used to test the performance of CRC\_RLS and its competing methods, including SRC [8], SVM, LRC (linear regression classification) [27] and NN. Note that LRC is actually an NS based method.

### 5.1. The role of sparsity: $l_1$ or $l_2$ ?

In this section, we study the role of sparsity in FR. Two representative face databases, Extended Yale B [23][24] and AR [25], are used (the experimental settings are described in Section 5.3). We use Eigenfaces of dimensionality 300 as the input facial features, and use all the training samples as the dictionary.



**Figure 4:** The recognition rates of SRC ( $l_1$ -regularized minimization) and CRC\_RLS ( $l_2$ -regularized minimization) versus the different values of  $\lambda$  on the (a) AR and (b) Extended Yale B databases; (c) the coding coefficients of a query sample.

The sparse coding of SRC in Eq. (1) can be equivalently written as  $(\hat{\boldsymbol{\alpha}}) = \text{arg min}_{\boldsymbol{\alpha}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right\}$ . We test the performance of SRC ( $l_1$ -regularized minimization) and CRC\_RLS ( $l_2$ -regularized minimization) by increasing the value of regularization parameter  $\lambda$ . The results on the AR and Extended Yale B databases are shown in Fig. 4(a) and Fig. 4(b), respectively. We can see that when  $\lambda=0$ , SRC and CRC\_RLS fail. When  $\lambda$  is assigned a small positive value, e.g., from 0.000001 to 0.1, good results can be achieved by SRC and CRC\_RLS. When  $\lambda$  is too big (e.g.,  $>0.1$ ) the recognition rates of both methods drop.

From Fig. 4 we can have the following findings. First, with the increase of sparsity ( $>0.000001$ ), no much benefit on recognition rate can be gained. Second,  $l_2$ -regularized minimization (i.e., CRC\_RLS) could get higher recognition rates than  $l_1$ -regularized minimization (i.e., SRC) in a broad range of  $\lambda$ . This implies that  $l_1$ -norm does not play the key role in face classification.

Fig. 4(c) plots the query sample’s coding coefficients by SRC and CRC\_RLS when they achieve their best results in the AR database. It can be seen that CRC\_RLS has much

weaker sparsity than SRC; however, it achieves not worse results. Again, sparsity of the representation coefficients is useful but not that crucial for FR. *What really crucial is the CR mechanism in both CRC\_RLS and SRC.*

## 5.2. Gender classification

In this section, we validate our claim in Section 3.1 that when the samples in each class are enough, there is no need to code the testing sample over the whole dictionary. We chose a non-occluded subset (14 images per subject) of AR [25] consisting of 50 male and 50 female subjects. Images of the first 25 males and 25 females were used for training, and the remaining images for testing. We used PCA to reduce the dimension of each image to 300. For this 2-class classification problem with enough training samples, we code the testing sample by each class' dictionary, and then classify it based on both the representation error and coefficient sparsity. That is, the query sample  $\mathbf{y}$  is classified to the class which gives the minimal  $r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{X}_i \boldsymbol{\alpha}\|_2 + \lambda \|\boldsymbol{\alpha}\|_1$  or  $r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{X}_i \boldsymbol{\alpha}\|_2 + \lambda \|\boldsymbol{\alpha}\|_2$ . The methods are then called  $L_1R$  (for  $l_1$ -regularized minimization) and  $L_2R$  (for  $l_2$ -regularized minimization).

We compare  $L_1R$  and  $L_2R$  with the CRC\_RLS, SRC, SVM, LRC and NN, and the results are listed in Table 3. One can see that  $L_1R$  and  $L_2R$  get the best results, which validates that coding on each class' dictionary is more powerful than coding on the whole dictionary when the training samples of each class are enough, no matter  $l_1$ - or  $l_2$ -regularized minimization is used. CRC\_RLS gets the second best result, about 1.4% higher than SRC.

**Table 3:** The results of different methods on gender classification using the AR database.

$L_1R$	$L_2R$	CRC_RLS	SRC	SVM	LRC	NN
<b>94.9%</b>	<b>94.9%</b>	<b>93.7%</b>	92.3%	92.4%	27.3%	90.7%

## 5.3. Face recognition

The proposed CRC\_RLS is then tested for FR. The Eigenface is used as the face feature.

*a) Extended Yale B Database:* The Extended Yale B [23] [24] database contains about 2,414 frontal face images of 38 individuals. We used the cropped and normalized face images of size 54×48, which were taken under varying illumination conditions. We randomly split the database into two halves. One half, which contains 32 images for each person, was used as the dictionary, and the other half was used for testing. Table 4 shows the recognition rates versus feature dimension by NN, LRC, SVM, SRC and CRC\_RLS. It can be seen that CRC\_RLS and SRC achieve very similar results in all dimensions (the difference of recognition rate is less than 0.5%). Since there are

relatively enough (32 per class) training samples, all the methods have not bad recognition rates.

**Table 4:** The face recognition results of different methods on the Extended Yale B database.

Dim	84	150	300
NN	85.8%	90.0%	91.6%
LRC	94.5%	95.1%	95.9%
SVM	94.9%	96.4%	97.0%
SRC	<b>95.5%</b>	<b>96.8%</b>	<b>97.9%</b>
CRC_RLS	95.0%	96.3%	<b>97.9%</b>

*2) AR database:* As in [8], a subset (with only illumination and expression changes) that contains 50 male subjects and 50 female subjects was chosen from the AR dataset [25] in our experiments. For each subject, the seven images from Session 1 were used for training, with the other seven images from Session 2 for testing. The images were cropped to 60×43. The comparison of competing methods is given in Table 5. We can see that CRC\_RLS achieves the best result when the dimensionality is 120 or 300. The recognition rates of CRC\_RLS and SRC are both at least 10% higher than other methods. This shows that CR does have much contribution to face classification.

**Table 5:** The face recognition results of different methods on the AR database.

Dim	54	120	300
NN	68.0%	70.1%	71.3%
LRC	71.0%	75.4%	76.0%
SVM	69.4%	74.5%	75.4%
SRC	<b>83.3%</b>	89.5%	93.3%
CRC_RLS	80.5%	<b>90.0%</b>	<b>93.7%</b>

**Table 6:** The face recognition results of different methods on the MPIE database.

	NN	LRC	SVM	SRC	CRC_RLS
S2	86.4%	87.1%	85.2%	93.9%	<b>94.1%</b>
S3	78.8%	81.9%	78.1%	<b>90.0%</b>	89.3%
S4	82.3%	84.3%	82.1%	<b>94.0%</b>	93.3%

*3) Multi PIE database:* The CMU Multi-PIE database [26] contains images of 337 subjects captured in four sessions with simultaneous variations in pose, expression, and illumination. In the experiments, all the 249 subjects in Session 1 were used. For the training set, we used the 14 frontal images with 14 illuminations<sup>2</sup> and neutral expression. For the testing sets, 10 typical frontal images<sup>3</sup> of illuminations taken with neutral expressions from Session 2 to Session 4 were used. The dimensionality of Eigenface is 300. Table 6 lists the recognition rates in three tests by the competing methods. The results validate that CRC\_RLS and SRC are the best in accuracy, with at least

<sup>2</sup> Illuminations {0,1,3,4,6,7,8,11,13,14,16,17,18,19}.

<sup>3</sup> Illuminations {0,2,4,6,8,10,12,14,16,18}.

6% improvement than the other three methods.

4) *FR with real face disguise*: As in [8], a subset from the AR database consisting of 1400 images from 100 subjects, 50 male and 50 female, is used here. 800 images (about 8 samples per subject) of non-occluded frontal views with various facial expressions were used for training, while the others with sunglasses and scarves (as shown in Fig. 5) were used for testing. The images were resized to 83×60. To handle the occlusion, SRC uses  $l_1$ -norm to fit the coding error and the sparse coding model is:  $(\hat{\alpha}) = \arg \min_{\alpha} \|\alpha\|_1$  s.t.  $\|y - X\alpha\|_1 < \varepsilon$  [8]. Note that the use of  $l_1$ -norm on the coding error increases much the complexity of SRC.

The results are shown in Table 7. Although CRC\_RLS is directly applied to the disguise face images, it gets the best result of FR with scarf disguise, outperforming SRC by a margin of 31%. For the case of FR with sunglasses, CRC\_RLS is worse than SRC, but still better than SVM. We also partitioned the face image into 8 sub-regions for testing (the partition is the same as that in [8]). Then both the recognition rates of CRC\_RLS and SRC are greater than 91%. These FR experiments with disguise again validate that CRC\_RLS is very competitive.



**Figure 5:** The testing samples with sunglasses and scarves in the AR database.

**Table 7:** The results of face recognition with real disguise using the AR database.

	Sunglass	Scarf
SVM	66.5%	16.5%
SRC	<b>87.0%</b>	59.5%
SRC (partitioned)	<b>97.5%</b>	93.5%
CRC_RLS	68.5%	<b>90.5%</b>
CRC_RLS (partitioned)	91.5%	<b>95%</b>

In all the above FR experiments, both CRC\_RLS and SRC are better than NN and LRC because of the benefit brought by CR. On the other hand, the result of CRC\_RLS is comparable to SRC, showing that the  $l_1$ -norm regularization does not bring more benefit than the simple  $l_2$ -norm regularization in FR.

## 5.4. Running time

At last, let's compare the running time of CRC\_RLS and SRC with various fast  $l_1$ -minimization methods, including  $l_1$ -ls [22], ALM [20], FISTA [19] and Homotopy[21]. We fix the dimensionality of Eigenface as 300. The recognition rates and speed of SRC and CRC\_RLS are listed in Table 8

(Extended Yale B), Table 9 (AR) and Table 10 (Multi-PIE), respectively. Note that the results in Table 10 are the averaged values of Sessions 2, 3 and 4.

**Table 8:** Recognition rate and speed on the Extended Yale B database.

	Recognition rate	Time
SRC( $l_1$ -ls)	<b>0.979</b>	5.3988 s
SRC(ALM)	<b>0.979</b>	0.128 s
SRC(FISTA)	0.914	0.1567 s
SRC(Homotopy)	0.945	0.0279 s
<b>CRC_RLS</b>	<b>0.979</b>	<b>0.0033 s</b>
<b>Speed-up</b>	<b>8.5 ~ 1636 times</b>	

**Table 9:** Recognition rate and speed on the AR database.

	Recognition rate	Time
SRC( $l_1$ -ls)	0.933	1.7878 s
SRC(ALM)	0.933	0.0578 s
SRC(FISTA)	0.6824	0.0457 s
SRC(Homotopy)	0.8212	0.0305 s
<b>CRC_RLS</b>	<b>0.937</b>	<b>0.0024 s</b>
<b>Speed-up</b>	<b>12.6 ~ 744.9 times</b>	

**Table 10:** Recognition rate and speed on the MPIE database.

	Recognition rate	Time
SRC( $l_1$ -ls)	<b>0.926</b>	21.2897 s
SRC(ALM)	0.9195	1.76 s
SRC(FISTA)	0.7955	1.636 s
SRC(Homotopy)	0.9017	0.5277 s
<b>CRC_RLS</b>	<b>0.922</b>	<b>0.0133 s</b>
<b>Speed-up</b>	<b>39.7 ~ 1600.7 times</b>	

On Yale B, CRC\_RLS, SRC( $l_1$ -ls) and SRC(ALM) achieve the best recognition rate (97.9%), but the speed of CRC\_RLS is 1636 and 38.8 times faster than SRC( $l_1$ -ls) and SRC(ALM). For the experiments on AR, CRC\_RLS has the best recognition rate and speed. SRC( $l_1$ -ls) is the second best but with the slowest speed. SRC(FISTA) and SRC(Homotopy) are much faster than SRC( $l_1$ -ls) but they have lower recognition rates. On Multi-PIE, CRC\_RLS achieves the second highest recognition rate (only 0.4% lower than SRC( $l_1$ -ls)) but it is significantly (more than 1600 times) faster than SRC( $l_1$ -ls). In this large-scale database, CRC\_RLS is about 40 times faster than SRC with the fastest implementation (i.e., Homotopy) with more than 2% improvement in recognition rate.

From the results in the above three tests, we can see that the speed-up of CRC\_RLS is more obvious as the scale (i.e., the number of classes or training samples) of face database increases. This implies that CRC\_RLS is more advantageous in practical large-scale FR applications.

## 6. Conclusion and discussions

This paper revealed that it is the collaborative representation (CR) mechanism, but not the  $l_1$ -norm

sparsity constraint, that truly improves the face recognition (FR) accuracy. We then presented a very simple yet very effective FR scheme, namely CR based classification with regularized least square (CRC\_RLS). Compared with the  $l_1$ -regularized sparse representation based classification (SRC), the  $l_2$ -regularized CRC\_RLS has very competitive FR accuracy but with significantly lower complexity. The extensive experimental results clearly demonstrated that CRC\_RLS is up to 1600 times faster than SRC without sacrificing recognition rate.

Apart from FR, our experiments on other types of signals (e.g., the human mouth odor signal classification for medical diagnosis) also showed that CRC or SRC works well. Statistically speaking, the norm (e.g.,  $l_1$  or  $l_2$ ) imposed on the coding coefficient and coding error depends on the distributions of them (e.g., Laplacian or Gaussian). Nonetheless, more investigations are to be made to further study the CRC scheme for various pattern classification problems, and this is one of our main objectives in the future work.

## References

- [1] B. Olshausen and D. Field. Sparse coding with an overcomplete basis set: a strategy employed by V1? *Vision Research*, 37(23):3311–3325, 1997.
- [2] M. Aharon, M. Elad, and A. M. Bruckstein. The K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representation. *IEEE SP*, 54(11):4311–4322, 2006.
- [3] J. Mairal, F. Bach, J. Ponce, G. Sapiro and A. Zisserman, Non-local sparse models for image restoration. In *ICCV* 2009.
- [4] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society B*, 58(1):267–288, 1996.
- [5] D. Donoho. For Most Large Underdetermined Systems of Linear Equations the Minimal  $l_1$ -Norm Solution is also the Sparsest Solution. *Comm. Pure and Applied Math.*, 59(6):797–829, 2006.
- [6] J. A. Tropp and S. J. Wright. Computational methods for sparse solution of linear inverse problems. *Proceedings of IEEE, Special Issue on Applications of Compressive Sensing & Sparse Representation*, 98(6):948–958, 2010.
- [7] K. Huang and S. Aviyente. Sparse representation for signal classification. In *NIPS*, 2006.
- [8] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE PAMI*, 31(2):210–227, 2009.
- [9] S. H. Gao, I. W-H. Tsang, and L-T. Chia. Kernel Sparse Representation for Image Classification and Face Recognition. In *ECCV*, 2010.
- [10] M. Yang and L. Zhang. Gabor Feature based Sparse Representation for Face Recognition with Gabor Occlusion Dictionary. In *ECCV*, 2010.
- [11] B. Cheng, J. Yang, S. Yan, Y. Fu, and T. Huang. Learning with  $l_1$ -graph for image analysis. *IEEE IP*, 19(4):858–866, 2010.
- [12] J. Yang, K. Yu, Y. Gong and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *CVPR* 2009.
- [13] J. Wright, Y. Ma, J. Mairal, G. Sapiro, T. Huang, and S. Yan. Sparse representation for computer vision and pattern recognition. *Proceedings of IEEE, Special Issue on Applications of Compressive Sensing & Sparse Representation*, 98(6):1031–1044, 2010.
- [14] J. Z. Huang, X. L. Huang, and D. Metaxas. Simultaneous image transformation and sparse representation recovery. In *CVPR* 2008.
- [15] A. Wagner, J. Wright, A. Ganesh, Z.H. Zhou, and Y. Ma, Towards a practical face recognition system: robust registration and illumination by sparse representation. In *CVPR* 2009.
- [16] Y. Peng, A. Ganesh, J. Wright, W. Xu, and Y. Ma. RASL: robust alignment by sparse and low-rank decomposition for linearly correlated images. Submitted to *IEEE PAMI*, 2010.
- [17] A. Y. Yang, A. Ganesh, Z. H. Zhou, S. S. Sastry, and Y. Ma. Fast  $l_1$ -minimization algorithms and application in robust face recognition, UC Berkeley, Tech. Rep.
- [18] S. J. Wright, R. D. Nowak, M. A. T. Figueiredo. Sparse reconstruction by separable approximation. In *ICASSP*, 2008.
- [19] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM. J. Imaging Science*, 2(1):183–202, 2009.
- [20] J. Yang and Y. Zhang. Alternating direction algorithms for  $l_1$ -problems in compressive sensing. (*preprint*) *arXiv:0912.1185*, 2009.
- [21] D. Malioutov, M. Cetin, and A. Willsky. Homotopy continuation for sparse signal representation. In *ICASSP*, 2005.
- [22] 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*, 1(4):606–617, 2007.
- [23] A. Georghiadis, P. Belhumeur, and D. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE PAMI*, 23(6):643–660, 2001.
- [24] K. Lee, J. Ho, and D. Kriegman. Acquiring linear subspaces for face recognition under variable lighting. *IEEE PAMI*, 27(5):684–698, 2005.
- [25] A. Martinez, and R. benavente. The AR face database. *CVC Tech. Report No. 24*, 1998.
- [26] R. Gross, I. Matthews. J. Cohn, T. Kanade, and S. Baker. Multi-PIE. *Image and Vision Computing*, 28:807–813, 2010.
- [27] I. Naseem, R. Togneri, and M. Bennamoun. Linear regression for face recognition. *IEEE PAMI*, 32(11):2106–2112, 2010.
- [28] G. H. Golub, and C. F. Van Loan, *Matrix Computation*, Johns Hopkins University Press, 1996.
- [29] R. Rigamonti, M. Brown and V. Lepetit. Are Sparse Representations Really Relevant for Image Classification? In *CVPR* 2011.
- [30] Q. Shi, A. Eriksson, A. Hengel, C. Shen. Is face recognition really a compressive sensing problem? In *CVPR* 2011.