

# Robust Face Recognition via Sparse Representation

John Wright, *Student Member, IEEE*, Allen Y. Yang, *Member, IEEE*,  
 Arvind Ganesh, *Student Member, IEEE*, S. Shankar Sastry, *Fellow, IEEE*, and  
 Yi Ma, *Senior Member, IEEE*

**Abstract**—We consider the problem of automatically recognizing human faces from frontal views with varying expression and illumination, as well as occlusion and disguise. We cast the recognition problem as one of classifying among multiple linear regression models and argue that new theory from sparse signal representation offers the key to addressing this problem. Based on a sparse representation computed by  $\ell^1$ -minimization, we propose a general classification algorithm for (image-based) object recognition. This new framework provides new insights into two crucial issues in face recognition: *feature extraction* and *robustness to occlusion*. For feature extraction, we show that if sparsity in the recognition problem is properly harnessed, the choice of features is no longer critical. What is critical, however, is whether the number of features is sufficiently large and whether the sparse representation is correctly computed. Unconventional features such as downsampled images and random projections perform just as well as conventional features such as Eigenfaces and Laplacianfaces, as long as the dimension of the feature space surpasses certain threshold, predicted by the theory of sparse representation. This framework can handle errors due to occlusion and corruption uniformly by exploiting the fact that these errors are often sparse with respect to the standard (pixel) basis. The theory of sparse representation helps predict how much occlusion the recognition algorithm can handle and how to choose the training images to maximize robustness to occlusion. We conduct extensive experiments on publicly available databases to verify the efficacy of the proposed algorithm and corroborate the above claims.

**Index Terms**—Face recognition, feature extraction, occlusion and corruption, sparse representation, compressed sensing,  $\ell^1$ -minimization, validation and outlier rejection.

## 1 INTRODUCTION

PARSIMONY has a rich history as a guiding principle for inference. One of its most celebrated instantiations, the principle of minimum description length in model selection [1], [2], stipulates that within a hierarchy of model classes, the model that yields the most compact representation should be preferred for decision-making tasks such as classification. A related, but simpler, measure of parsimony in high-dimensional data processing seeks models that depend on only a few of the observations, selecting a small subset of features for classification or visualization (e.g., Sparse PCA [3], [4] among others). Such sparse feature selection methods are, in a sense, dual to the support vector machine (SVM) approach in [5] and [6], which instead selects a small subset of relevant training examples to characterize the decision boundary between classes. While these works comprise only a small fraction of the literature on parsimony for inference, they do serve to illustrate a common theme: all of them use parsimony as a principle for

choosing a limited subset of features or models from the training data, rather than directly using the data for representing or classifying an input (test) signal.

The role of parsimony in human perception has also been strongly supported by studies of human vision. Investigators have recently revealed that in both low-level and midlevel human vision [7], [8], many neurons in the visual pathway are selective for a variety of specific stimuli, such as color, texture, orientation, scale, and even view-tuned object images. Considering these neurons to form an overcomplete dictionary of base signal elements at each visual stage, the firing of the neurons with respect to a given input image is typically highly sparse.

In the statistical signal processing community, the algorithmic problem of computing sparse linear representations with respect to an overcomplete dictionary of base elements or signal atoms has seen a recent surge of interest [9], [10], [11], [12].<sup>1</sup> Much of this excitement centers around the discovery that whenever the optimal representation is sufficiently sparse, it can be efficiently computed by convex optimization [9], even though this problem can be extremely difficult in the general case [13]. The resulting optimization problem, similar to the Lasso in statistics

- J. Wright, A. Ganesh, and Y. Ma are with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, 1308 West Main Street, Urbana, IL 61801. E-mail: {jwright, abalasu2, yima}@uiuc.edu.
- A. Yang and S. Sastry are with the Department of Electrical Engineering and Computer Science, University of California, Berkeley, Berkeley, CA 94720. e-mail: {yang, sastry}@eecs.berkeley.edu.

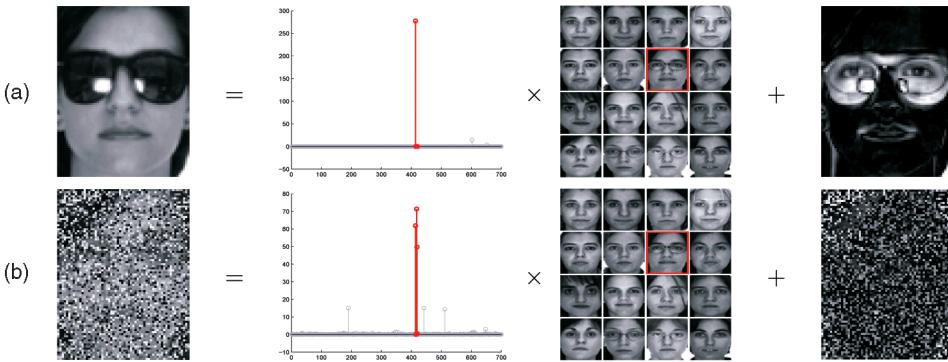
Manuscript received 13 Aug. 2007; revised 18 Jan. 2008; accepted 20 Mar. 2008; published online 26 Mar. 2008.

Recommended for acceptance by M.-H. Yang.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2007-08-0500.

Digital Object Identifier no. 10.1109/TPAMI.2008.79.

1. In the literature, the terms “sparse” and “representation” have been used to refer to a number of similar concepts. Throughout this paper, we will use the term “sparse representation” to refer specifically to an expression of the input signal as a linear combination of base elements in which many of the coefficients are zero. In most cases considered, the percentage of nonzero coefficients will vary between zero and  $\approx 30$  percent. However, in characterizing the breakdown point of our algorithms, we will encounter cases with up to 70 percent nonzeros.



**Fig. 1. Overview of our approach.** Our method represents a test image (left), which is (a) potentially occluded or (b) corrupted, as a sparse linear combination of all the training images (middle) plus sparse errors (right) due to occlusion or corruption. Red (darker) coefficients correspond to training images of the correct individual. Our algorithm determines the true identity (indicated with a red box at second row and third column) from 700 training images of 100 individuals (7 each) in the standard AR face database.

[12], [14] penalizes the  $\ell^1$ -norm of the coefficients in the linear combination, rather than the directly penalizing the number of nonzero coefficients (i.e., the  $\ell^0$ -norm).

The original goal of these works was not inference or classification per se, but rather representation and compression of signals, potentially using lower sampling rates than the Shannon-Nyquist bound [15]. Algorithm performance was therefore measured in terms of sparsity of the representation and fidelity to the original signals. Furthermore, individual base elements in the dictionary were not assumed to have any particular semantic meaning—they are typically chosen from standard bases (e.g., Fourier, Wavelet, Curvelet, and Gabor), or even generated from random matrices [11], [15]. Nevertheless, the sparsest representation is naturally discriminative: among all subsets of base vectors, it selects the subset which most compactly expresses the input signal and rejects all other possible but less compact representations.

In this paper, we exploit the discriminative nature of sparse representation to **perform classification**. Instead of using the generic dictionaries discussed above, we represent the test sample in an overcomplete dictionary whose base elements are *the training samples themselves*. If sufficient training samples are available from each class,<sup>2</sup> it will be possible to represent the test samples as a linear combination of just those training samples from the same class. This representation is naturally sparse, involving only a small fraction of the overall training database. We argue that in many problems of interest, it is actually the *sparsest* linear representation of the test sample in terms of this dictionary and can be recovered efficiently via  $\ell^1$ -minimization. Seeking the sparsest representation therefore automatically discriminates between the various classes present in the training set. Fig. 1 illustrates this simple idea using face recognition as an example. Sparse representation also provides a simple and surprisingly effective means of rejecting invalid test samples not arising from any class in the training database: these samples' sparsest representations tend to involve many dictionary elements, spanning multiple classes.

2. In contrast, methods such as that in [16] and [17] that utilize only a single training sample per class face a more difficult problem and generally incorporate more explicit prior knowledge about the types of variation that could occur in the test sample.

Our use of sparsity for classification differs significantly from the various parsimony principles discussed above. Instead of using sparsity to identify a relevant model or relevant features that can later be used for classifying *all* test samples, it uses the sparse representation of each individual test sample directly for classification, adaptively selecting the training samples that give the most compact representation. The proposed classifier can be considered a generalization of popular classifiers such as *nearest neighbor* (NN) [18] and *nearest subspace* (NS) [19] (i.e., minimum distance to the subspace spanned all training samples from each object class). NN classifies the test sample based on the best representation in terms of a single training sample, whereas NS classifies based on the best linear representation in terms of all the training samples in each class. The *nearest feature line* (NFL) algorithm [20] strikes a balance between these two extremes, classifying based on the best affine representation in terms of a pair of training samples. Our method strikes a similar balance but considers all possible supports (within each class or across multiple classes) and adaptively chooses the minimal number of training samples needed to represent each test sample.<sup>3</sup>

We will motivate and study this new approach to classification within the context of automatic face recognition. Human faces are arguably the most extensively studied object in image-based recognition. This is partly due to the remarkable face recognition capability of the human visual system [21] and partly due to numerous important applications for face recognition technology [22]. In addition, technical issues associated with face recognition are representative of object recognition and even data classification in general. Conversely, the theory of sparse representation and compressed sensing yields new insights into two crucial issues in automatic face recognition: the role of feature extraction and the difficulty due to occlusion.

**The role of feature extraction.** The question of *which low-dimensional features of an object image are the most relevant or informative for classification* is a central issue in face recognition and in object recognition in general. An enormous volume of literature has been devoted to investigate various data-dependent feature transformations for projecting the

3. The relationship between our method and NN, NS, and NFL is explored more thoroughly in the supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>.

high-dimensional test image into lower dimensional feature spaces: examples include Eigenfaces [23], Fisherfaces [24], Laplacianfaces [25], and a host of variants [26], [27]. With so many proposed features and so little consensus about which are better or worse, practitioners lack guidelines to decide which features to use. However, within our proposed framework, the theory of compressed sensing implies that *the precise choice of feature space is no longer critical*: Even random features contain enough information to recover the sparse representation and hence correctly classify any test image. What is critical is that the dimension of the feature space is sufficiently large and that the sparse representation is correctly computed.

**Robustness to occlusion.** Occlusion poses a significant obstacle to robust real-world face recognition [16], [28], [29]. This difficulty is mainly due to the unpredictable nature of the error incurred by occlusion: it may affect any part of the image and may be arbitrarily large in magnitude. Nevertheless, this error typically corrupts only a fraction of the image pixels and is therefore sparse in the standard basis given by individual pixels. When the error has such a sparse representation, it can be handled uniformly within our framework: the basis in which the error is sparse can be treated as a special class of training samples. The subsequent sparse representation of an occluded test image with respect to this expanded dictionary (training images plus error basis) naturally separates the component of the test image arising due to occlusion from the component arising from the identity of the test subject (see Fig. 1 for an example). In this context, the theory of sparse representation and compressed sensing characterizes when such *source-and-error separation* can take place and therefore how much occlusion the resulting recognition algorithm can tolerate.

**Organization of this paper.** In Section 2, we introduce a basic general framework for classification using sparse representation, applicable to a wide variety of problems in image-based object recognition. We will discuss why the sparse representation can be computed by  $\ell^1$ -minimization and how it can be used for classifying and validating any given test sample. Section 3 shows how to apply this general classification framework to study two important issues in image-based face recognition: feature extraction and robustness to occlusion. In Section 4, we verify the proposed method with extensive experiments on popular face data sets and comparisons with many other state-of-the-art face recognition techniques. Further connections between our method, NN, and NS are discussed in the supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>.

While the proposed method is of broad interest to object recognition in general, the studies and experimental results in this paper are confined to human frontal face recognition. We will deal with illumination and expressions, but we do not explicitly account for object pose nor rely on any 3D model of the face. The proposed algorithm is robust to small variations in pose and displacement, for example, due to registration errors. However, we do assume that detection, cropping, and normalization of the face have been performed prior to applying our algorithm.

## 2 CLASSIFICATION BASED ON SPARSE REPRESENTATION

A basic problem in object recognition is to use labeled training samples from  $k$  distinct object classes to correctly determine the class to which a new test sample belongs. We arrange the given  $n_i$  training samples from the  $i$ th class as columns of a matrix  $A_i = [v_{i,1}, v_{i,2}, \dots, v_{i,n_i}] \in \mathbb{R}^{m \times n_i}$ . In the context of face recognition, we will identify a  $w \times h$  grayscale image with the vector  $v \in \mathbb{R}^m$  ( $m = wh$ ) given by stacking its columns; the columns of  $A_i$  are then the training face images of the  $i$ th subject.

### 2.1 Test Sample as a Sparse Linear Combination of Training Samples

An immense variety of statistical, generative, or discriminative models have been proposed for exploiting the structure of the  $A_i$  for recognition. One particularly simple and effective approach models the samples from a single class as lying on a linear subspace. Subspace models are flexible enough to capture much of the variation in real data sets and are especially well motivated in the context of face recognition, where it has been observed that the images of faces under varying lighting and expression lie on a special low-dimensional subspace [24], [30], often called a *face subspace*. Although the proposed framework and algorithm can also apply to multimodal or nonlinear distributions (see the supplementary appendix for more detail, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>), for ease of presentation, we shall first assume that the training samples from a single class do lie on a subspace. This is the only prior knowledge about the training samples we will be using in our solution.<sup>4</sup>

Given sufficient training samples of the  $i$ th object class,  $A_i = [v_{i,1}, v_{i,2}, \dots, v_{i,n_i}] \in \mathbb{R}^{m \times n_i}$ , any new (test) sample  $y \in \mathbb{R}^m$  from the same class will approximately lie in the linear span of the training samples<sup>5</sup> associated with object  $i$ :

$$y = \alpha_{i,1}v_{i,1} + \alpha_{i,2}v_{i,2} + \cdots + \alpha_{i,n_i}v_{i,n_i}, \quad (1)$$

for some scalars,  $\alpha_{i,j} \in \mathbb{R}$ ,  $j = 1, 2, \dots, n_i$ .

Since the membership  $i$  of the test sample is initially unknown, we define a new matrix  $A$  for the entire training set as the concatenation of the  $n$  training samples of all  $k$  object classes:

$$A \doteq [A_1, A_2, \dots, A_k] = [v_{1,1}, v_{1,2}, \dots, v_{k,n_k}]. \quad (2)$$

Then, the linear representation of  $y$  can be rewritten in terms of all training samples as

$$y = Ax_0 \in \mathbb{R}^m, \quad (3)$$

where  $x_0 = [0, \dots, 0, \alpha_{i,1}, \alpha_{i,2}, \dots, \alpha_{i,n_i}, 0, \dots, 0]^T \in \mathbb{R}^n$  is a coefficient vector whose entries are zero except those associated with the  $i$ th class.

4. In face recognition, we actually do not need to know whether the linear structure is due to varying illumination or expression, since we do not rely on domain-specific knowledge such as an illumination model [31] to eliminate the variability in the training and testing images.

5. One may refer to [32] for how to choose the training images to ensure this property for face recognition. Here, we assume that such a training set is given.

As the entries of the vector  $x_0$  encode the identity of the test sample  $y$ , it is tempting to attempt to obtain it by solving the linear system of equations  $y = Ax$ . Notice, though, that using the entire training set to solve for  $x$  represents a significant departure from one sample or one class at a time methods such as NN and NS. We will later argue that one can obtain a more discriminative classifier from such a global representation. We will demonstrate its superiority over these local methods (NN or NS) both for identifying objects represented in the training set and for rejecting outlying samples that do not arise from any of the classes present in the training set. These advantages can come without an increase in the order of growth of the computation. As we will see, the complexity remains linear in the size of training set.

Obviously, if  $m > n$ , the system of equations  $y = Ax$  is overdetermined, and the correct  $x_0$  can usually be found as its unique solution. We will see in Section 3, however, that in robust face recognition, the system  $y = Ax$  is typically underdetermined, and so, its solution is not unique.<sup>6</sup> Conventionally, this difficulty is resolved by choosing the minimum  $\ell^2$ -norm solution:

$$(\ell^2) : \hat{x}_2 = \arg \min \|x\|_2 \text{ subject to } Ax = y. \quad (4)$$

While this optimization problem can be easily solved (via the pseudoinverse of  $A$ ), the solution  $\hat{x}_2$  is not especially informative for recognizing the test sample  $y$ . As shown in Example 1,  $\hat{x}_2$  is generally *dense*, with large nonzero entries corresponding to training samples from many different classes. To resolve this difficulty, we instead exploit the following simple observation: A valid test sample  $y$  can be sufficiently represented using only the training samples from the same class. This representation is naturally *sparse* if the number of object classes  $k$  is reasonably large. For instance, if  $k = 20$ , only 5 percent of the entries of the desired  $x_0$  should be nonzero. The more sparse the recovered  $x_0$  is, the easier will it be to accurately determine the identity of the test sample  $y$ .<sup>7</sup>

This motivates us to seek the sparsest solution to  $y = Ax$ , solving the following optimization problem:

$$(\ell^0) : \hat{x}_0 = \arg \min \|x\|_0 \text{ subject to } Ax = y, \quad (5)$$

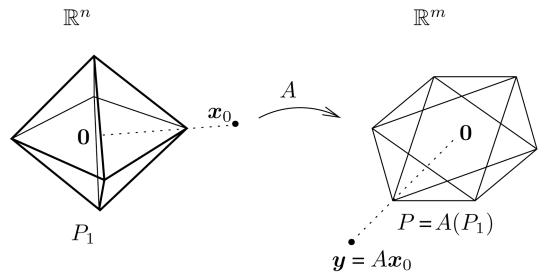
where  $\|\cdot\|_0$  denotes the  $\ell^0$ -norm, which counts the number of nonzero entries in a vector. In fact, if the columns of  $A$  are in general position, then whenever  $y = Ax$  for some  $x$  with less than  $m/2$  nonzeros,  $x$  is the unique sparsest solution:  $\hat{x}_0 = x$  [33]. However, the problem of finding the sparsest solution of an underdetermined system of linear equations is NP-hard and difficult even to approximate [13]: that is, in the general case, no known procedure for finding the sparsest solution is significantly more efficient than exhausting all subsets of the entries for  $x$ .

## 2.2 Sparse Solution via $\ell^1$ -Minimization

Recent development in the emerging theory of *sparse representation and compressed sensing* [9], [10], [11] reveals

6. Furthermore, even in the overdetermined case, such a linear equation may not be perfectly satisfied in the presence of data noise (see Section 2.2.2).

7. This intuition holds only when the size of the database is fixed. For example, if we are allowed to append additional irrelevant columns to  $A$ , we can make the solution  $x_0$  have a smaller fraction of nonzeros, but this does not make  $x_0$  more informative for recognition.



**Fig. 2. Geometry of sparse representation via  $\ell^1$ -minimization.** The  $\ell^1$ -minimization determines which facet (of the lowest dimension) of the polytope  $A(P_\alpha)$ , the point  $y/\|y\|_1$ , lies in. The test sample vector  $y$  is represented as a linear combination of just the vertices of that facet, with coefficients  $x_0$ .

that if the solution  $x_0$  sought is *sparse enough*, the solution of the  $\ell^1$ -minimization problem (5) is equal to the solution to the following  $\ell^1$ -minimization problem:

$$(\ell^1) : \hat{x}_1 = \arg \min \|x\|_1 \text{ subject to } Ax = y. \quad (6)$$

This problem can be solved in polynomial time by standard linear programming methods [34]. Even more efficient methods are available when the solution is known to be very sparse. For example, homotopy algorithms recover solutions with  $t$  nonzeros in  $O(t^3 + n)$  time, linear in the size of the training set [35].

### 2.2.1 Geometric Interpretation

Fig. 2 gives a geometric interpretation (essentially due to [36]) of why minimizing the  $\ell^1$ -norm correctly recovers sufficiently sparse solutions. Let  $P_\alpha$  denote the  $\ell^1$ -ball (or crosspolytope) of radius  $\alpha$ :

$$P_\alpha := \{x : \|x\|_1 \leq \alpha\} \subset \mathbb{R}^n. \quad (7)$$

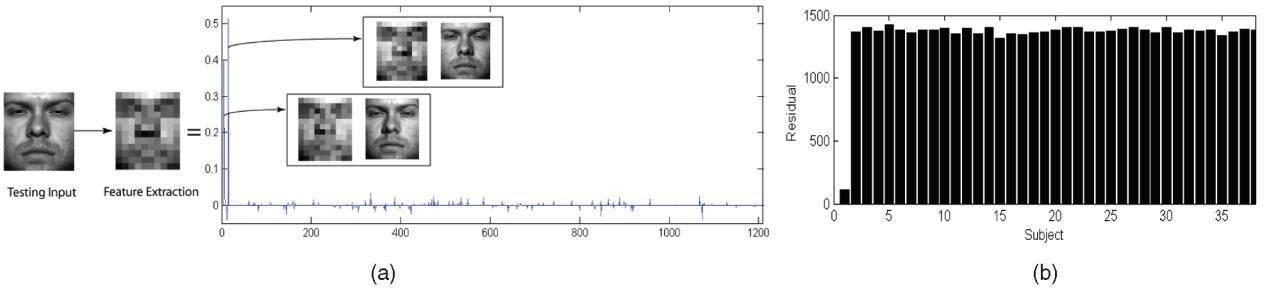
In Fig. 2, the unit  $\ell^1$ -ball  $P_1$  is mapped to the polytope  $P := A(P_1) \subset \mathbb{R}^m$ , consisting of all  $y$  that satisfy  $y = Ax$  for some  $x$  whose  $\ell^1$ -norm is  $\leq 1$ .

The geometric relationship between  $P_\alpha$  and the polytope  $A(P_\alpha)$  is invariant to scaling. That is, if we scale  $P_\alpha$ , its image under multiplication by  $A$  is also scaled by the same amount. Geometrically, finding the minimum  $\ell^1$ -norm solution  $\hat{x}_1$  to (6) is equivalent to expanding the  $\ell^1$ -ball  $P_\alpha$  until the polytope  $A(P_\alpha)$  first touches  $y$ . The value of  $\alpha$  at which this occurs is exactly  $\|\hat{x}_1\|_1$ .

Now, suppose that  $y = Ax_0$  for some sparse  $x_0$ . We wish to know when solving (6) correctly recovers  $x_0$ . This question is easily resolved from the geometry of that in Fig. 2: Since  $\hat{x}_1$  is found by expanding both  $P_\alpha$  and  $A(P_\alpha)$  until a point of  $A(P_\alpha)$  touches  $y$ , the  $\ell^1$ -minimizer  $\hat{x}_1$  must generate a point  $A\hat{x}_1$  on the boundary of  $P$ .

Thus,  $\hat{x}_1 = x_0$  if and only if the point  $A(x_0/\|x_0\|_1)$  lies on the boundary of the polytope  $P$ . For the example shown in Fig. 2, it is easy to see that the  $\ell^1$ -minimization recovers all  $x_0$  with only one nonzero entry. This equivalence holds because all of the vertices of  $P_1$  map to points on the boundary of  $P$ .

In general, if  $A$  maps all  $t$ -dimensional facets of  $P_1$  to facets of  $P$ , the polytope  $P$  is referred to as (*centrally*)  $t$ -neighborly [36]. From the above, we see that the  $\ell^1$ -minimization (6) correctly recovers all  $x_0$  with  $\leq t+1$  nonzeros if and only if  $P$  is  $t$ -neighborly, in which case, it is



**Fig. 3. A valid test image.** (a) Recognition with  $12 \times 10$  downsampled images as features. The test image  $y$  belongs to subject 1. The values of the sparse coefficients recovered from Algorithm 1 are plotted on the right together with the two training examples that correspond to the two largest sparse coefficients. (b) The residuals  $r_i(y)$  of a test image of subject 1 with respect to the projected sparse coefficients  $\delta_i(\hat{x})$  by  $\ell^1$ -minimization. The ratio between the two smallest residuals is about 1:8.6.

equivalent to the  $\ell^0$ -minimization (5).<sup>8</sup> This condition is surprisingly common: even polytopes  $P$  given by random matrices (e.g., uniform, Gaussian, and partial Fourier) are highly neighborly [15], allowing correct recover of sparse  $x_0$  by  $\ell^1$ -minimization.

Unfortunately, there is no known algorithm for efficiently verifying the neighborliness of a given polytope  $P$ . The best known algorithm is combinatorial, and therefore, only practical when the dimension  $m$  is moderate [37]. When  $m$  is large, it is known that with overwhelming probability, the neighborliness of a randomly chosen polytope  $P$  is loosely bounded between

$$c \cdot m < t < [(m+1)/3], \quad (8)$$

for some small constant  $c > 0$  (see [9] and [36]). Loosely speaking, as long as the number of nonzero entries of  $x_0$  is a small fraction of the dimension  $m$ ,  $\ell^1$ -minimization will recover  $x_0$ .

### 2.2.2 Dealing with Small Dense Noise

So far, we have assumed that (3) holds exactly. Since real data are noisy, it may not be possible to express the test sample exactly as a sparse superposition of the training samples. The model (3) can be modified to explicitly account for small possibly dense noise by writing

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{z}, \quad (9)$$

where  $\mathbf{z} \in \mathbb{R}^m$  is a noise term with bounded energy  $\|\mathbf{z}\|_2 < \varepsilon$ . The sparse solution  $\mathbf{x}_0$  can still be approximately recovered by solving the following stable  $\ell^1$ -minimization problem:

$$(\ell_s^1) : \quad \hat{\mathbf{x}}_1 = \arg \min \|\mathbf{x}\|_1 \quad \text{subject to} \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon. \quad (10)$$

This convex optimization problem can be efficiently solved via second-order cone programming [34] (see Section 4 for our algorithm of choice). The solution of  $(\ell_s^1)$  is guaranteed to approximately recovery sparse solutions in ensembles of random matrices  $A$  [38]: There are constants  $\rho$  and  $\zeta$  such that with overwhelming probability, if  $\|\mathbf{x}_0\|_0 \leq \rho m$  and  $\|\mathbf{z}\|_2 \leq \varepsilon$ , then the computed  $\hat{\mathbf{x}}_1$  satisfies

$$\|\hat{\mathbf{x}}_1 - \mathbf{x}_0\|_2 \leq \zeta \varepsilon. \quad (11)$$

8. Thus, neighborliness gives a necessary and sufficient condition for sparse recovery. The restricted isometry properties often used in analyzing the performance of  $\ell^1$ -minimization in random matrix ensembles (e.g., [15]) give sufficient, but *not* necessary, conditions.

### 2.3 Classification Based on Sparse Representation

Given a new test sample  $\mathbf{y}$  from one of the classes in the training set, we first compute its sparse representation  $\hat{\mathbf{x}}_1$  via (6) or (10). Ideally, the nonzero entries in the estimate  $\hat{\mathbf{x}}_1$  will all be associated with the columns of  $A$  from a single object class  $i$ , and we can easily assign the test sample  $\mathbf{y}$  to that class. However, noise and modeling error may lead to small nonzero entries associated with multiple object classes (see Fig. 3). Based on the global sparse representation, one can design many possible classifiers to resolve this. For instance, we can simply assign  $\mathbf{y}$  to the object class with the single largest entry in  $\hat{\mathbf{x}}_1$ . However, such heuristics do not harness the subspace structure associated with images in face recognition. To better harness such linear structure, we instead classify  $\mathbf{y}$  based on how well the coefficients associated with all training samples of each object reproduce  $\mathbf{y}$ .

For each class  $i$ , let  $\delta_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be the characteristic function that selects the coefficients associated with the  $i$ th class. For  $\mathbf{x} \in \mathbb{R}^n$ ,  $\delta_i(\mathbf{x}) \in \mathbb{R}^n$  is a new vector whose only nonzero entries are the entries in  $\mathbf{x}$  that are associated with class  $i$ . Using only the coefficients associated with the  $i$ th class, one can approximate the given test sample  $\mathbf{y}$  as  $\hat{\mathbf{y}}_i = \mathbf{A}\delta_i(\hat{\mathbf{x}}_1)$ . We then classify  $\mathbf{y}$  based on these approximations by assigning it to the object class that minimizes the residual between  $\mathbf{y}$  and  $\hat{\mathbf{y}}_i$ :

$$\min_i r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{A}\delta_i(\hat{\mathbf{x}}_1)\|_2, \quad (12)$$

Algorithm 1 below summarizes the complete recognition procedure. Our implementation minimizes the  $\ell^1$ -norm via a primal-dual algorithm for linear programming based on [39] and [40].

#### Algorithm 1. Sparse Representation-based Classification (SRC)

##### 1: Input: a matrix of training samples

$A = [A_1, A_2, \dots, A_k] \in \mathbb{R}^{m \times n}$  for  $k$  classes, a test sample  $\mathbf{y} \in \mathbb{R}^m$ , (and an optional error tolerance  $\varepsilon > 0$ .)

2: Normalize the columns of  $A$  to have unit  $\ell^2$ -norm.

3: Solve the  $\ell^1$ -minimization problem:

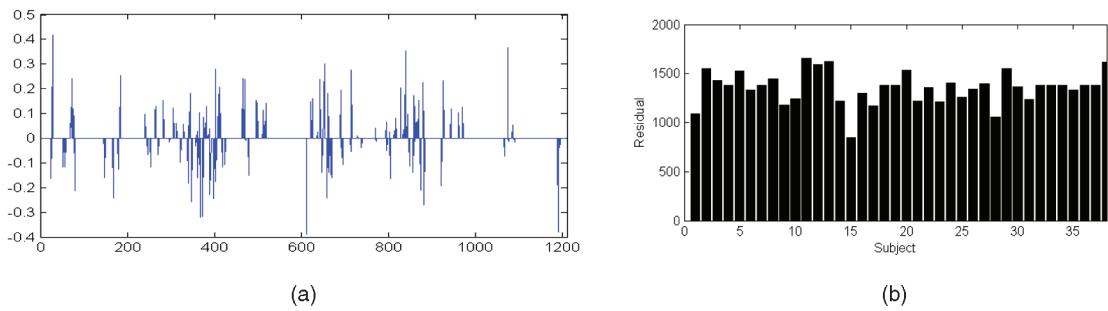
$$\hat{\mathbf{x}}_1 = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{y}. \quad (13)$$

(Or alternatively, solve

$$\hat{\mathbf{x}}_1 = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon.)$$

4: Compute the residuals  $r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{A}\delta_i(\hat{\mathbf{x}}_1)\|_2$  for  $i = 1, \dots, k$ .

5: Output: identity( $\mathbf{y}$ ) =  $\arg \min_i r_i(\mathbf{y})$ .



**Fig. 4. Nonsparsity of the  $\ell^2$ -minimizer.** (a) Coefficients from  $\ell^2$ -minimization using the same test image as Fig. 3. The recovered solution is not sparse and, hence, less informative for recognition (large coefficients do not correspond to training images of this test subject). (b) The residuals of the test image from subject 1 with respect to the projection  $\delta_i(\hat{x})$  of the coefficients obtained by  $\ell^2$ -minimization. The ratio between the two smallest residuals is about 1:1.3. The smallest residual is not associated with subject 1.

**Example 1 ( $\ell^1$ -minimization versus  $\ell^2$ -minimization).** To illustrate how Algorithm 1 works, we randomly select half of the 2,414 images in the Extended Yale B database as the training set and the rest for testing. In this example, we subsample the images from the original  $192 \times 168$  to size  $12 \times 10$ . The pixel values of the downsampled image are used as 120-D features—stacked as columns of the matrix  $A$  in the algorithm. Hence, matrix  $A$  has size  $120 \times 1,207$ , and the system  $y = Ax$  is underdetermined. Fig. 3a illustrates the sparse coefficients recovered by Algorithm 1 for a test image from the first subject. The figure also shows the features and the original images that correspond to the two largest coefficients. The two largest coefficients are both associated with training samples from subject 1. Fig. 3b shows the residuals with respect to the 38 projected coefficients  $\delta_i(\hat{x}_1)$ ,  $i = 1, 2, \dots, 38$ . With  $12 \times 10$  downsampled images as features, Algorithm 1 achieves an overall recognition rate of 92.1 percent across the Extended Yale B database. (See Section 4 for details and performance with other features such as Eigenfaces and Fisherfaces, as well as comparison with other methods.) Whereas the more conventional minimum  $\ell^2$ -norm solution to the underdetermined system  $y = Ax$  is typically quite dense, minimizing the  $\ell^1$ -norm favors sparse solutions and provably recovers the sparsest solution when this solution is sufficiently sparse. To illustrate this contrast, Fig. 4a shows the coefficients of the same test image given by the conventional  $\ell^2$ -minimization (4), and Fig. 4b shows the corresponding residuals with respect to the 38 subjects. The coefficients are much less sparse than those given by  $\ell^1$ -minimization (in Fig. 3), and the dominant coefficients are not associated with subject 1. As a result, the smallest residual in Fig. 4 does not correspond to the correct subject (subject 1).

#### 2.4 Validation Based on Sparse Representation

Before classifying a given test sample, we must first decide if it is a valid sample from one of the classes in the data set. The ability to detect and then reject invalid test samples, or “outliers,” is crucial for recognition systems to work in real-world situations. A face recognition system, for example, could be given a face image of a subject that is not in the database or an image that is not a face at all.

Systems based on conventional classifiers such as NN or NS, often use the residuals  $r_i(y)$  for validation, in addition to identification. That is, the algorithm accepts or rejects a

test sample based on how small the smallest residual is. However, each residual  $r_i(y)$  is computed without any knowledge of images of other object classes in the training data set and only measures similarity between the test sample and each individual class.

In the sparse representation paradigm, the coefficients  $\hat{x}$  are computed globally, in terms of images of all classes. In a sense, it can harness the joint distribution of all classes for validation. We contend that the coefficients  $\hat{x}$  are better statistics for validation than the residuals. Let us first see this through an example.

**Example 2 (concentration of sparse coefficients).** We randomly select an irrelevant image from Google and downsample it to  $12 \times 10$ . We then compute the sparse representation of the image against the same Extended Yale B training data, as in Example 1. Fig. 5a plots the obtained coefficients, and Fig. 5b plots the corresponding residuals. Compared to the coefficients of a valid test image in Fig. 3, notice that the coefficients  $\hat{x}$  here are not concentrated on any one subject and instead spread widely across the entire training set. Thus, the distribution of the estimated sparse coefficients  $\hat{x}$  contains important information about the validity of the test image: a valid test image should have a sparse representation whose nonzero entries concentrate mostly on one subject, whereas an invalid image has sparse coefficients spread widely among multiple subjects.

To quantify this observation, we define the following measure of how concentrated the coefficients are on a single class in the data set:

**Definition 1 (sparsity concentration index (SCI)).** The SCI of a coefficient vector  $x \in \mathbb{R}^n$  is defined as

$$\text{SCI}(x) = \frac{k \cdot \max_i \|\delta_i(x)\|_1 / \|x\|_1 - 1}{k - 1} \in [0, 1]. \quad (14)$$

For a solution  $\hat{x}$  found by Algorithm 1, if  $\text{SCI}(\hat{x}) = 1$ , the test image is represented using only images from a single object, and if  $\text{SCI}(\hat{x}) = 0$ , the sparse coefficients are spread evenly over all classes.<sup>9</sup> We choose a threshold  $\tau \in (0, 1)$  and accept a test image as valid if

<sup>9</sup> Directly choosing  $x$  to minimize the SCI might produce more concentrated coefficients; however, the SCI is highly nonconvex and difficult to optimize. For valid test images, minimizing the  $\ell^1$ -norm already produces representations that are well-concentrated on the correct subject class.

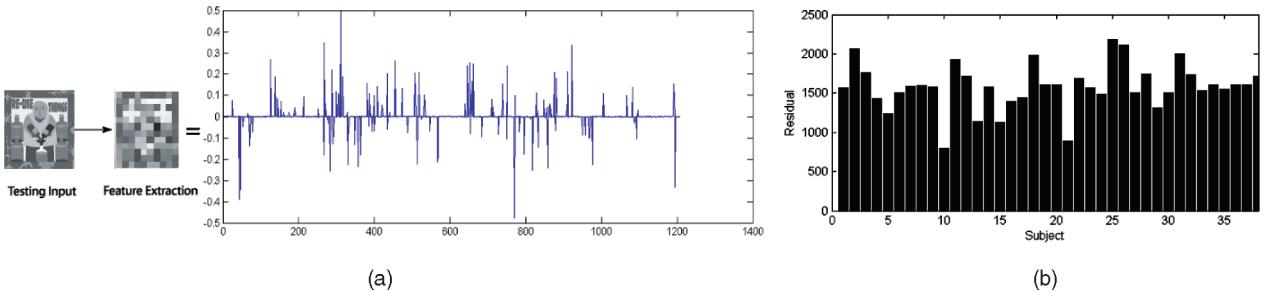


Fig. 5. **Example of an invalid test image.** (a) Sparse coefficients for the invalid test image with respect to the same training data set from Example 1. The test image is a randomly selected irrelevant image. (b) The residuals of the invalid test image with respect to the projection  $\delta_i(\hat{x})$  of the sparse representation computed by  $\ell^1$ -minimization. The ratio of the two smallest residuals is about 1:1.2.

$$\text{SCI}(\hat{x}) \geq \tau, \quad (15)$$

and otherwise reject as invalid. In step 5 of Algorithm 1, one may choose to output the identity of  $y$  only if it passes this criterion.

Unlike NN or NS, this new rule avoids the use of the residuals  $r_i(y)$  for validation. Notice that in Fig. 5, even for a nonface image, with a large training set, the smallest residual of the invalid test image is not so large. Rather than relying on a single statistic for both validation and identification, our approach separates the information required for these tasks: **the residuals for identification and the sparse coefficients for validation.**<sup>10</sup> In a sense, the residual measures how well the representation approximates the test image; and the sparsity concentration index measures how good the representation itself is, in terms of localization.

One benefit to this approach to validation is improved performance against generic objects that are similar to multiple object classes. For example, in face recognition, a generic face might be rather similar to some of the subjects in the data set and may have small residuals with respect to their training images. Using residuals for validation more likely leads to a false positive. However, a generic face is unlikely to pass the new validation rule as a good representation of it typically requires contribution from images of multiple subjects in the data set. Thus, the new rule can better judge whether the test image is a generic face or the face of one particular subject in the data set. In Section 4.7, we will demonstrate that the new validation rule outperforms the NN and NS methods, with as much as 10-20 percent improvement in verification rate for a given false accept rate (see Fig. 14 in Section 4 or Fig. 18 in the supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>).

### 3 TWO FUNDAMENTAL ISSUES IN FACE RECOGNITION

In this section, we study the implications of the above general classification framework for two critical issues in face recognition: 1) the choice of feature transformation, and 2) robustness to corruption, occlusion, and disguise.

10. We find empirically that this separation works well enough in our experiments with face images. However, it is possible that better validation and identification rules can be contrived from using the residual and the sparsity together.

### 3.1 The Role of Feature Extraction

In the computer vision literature, numerous feature extraction schemes have been investigated for finding projections that better separate the classes in lower dimensional spaces, which are often referred to as *feature spaces*. One class of methods extracts holistic face features such as Eigenfaces[23], Fisherfaces [24], and Laplacianfaces [25]. Another class of methods tries to extract meaningful partial facial features (e.g., patches around eyes or nose) [21], [41] (see Fig. 6 for some examples). Traditionally, when feature extraction is used in conjunction with simple classifiers such as NN and NS, the choice of feature transformation is considered critical to the success of the algorithm. This has led to the development of a wide variety of increasingly complex feature extraction methods, including nonlinear and kernel features [42], [43]. In this section, we reexamine the role of feature extraction within the new sparse representation framework for face recognition.

One benefit of feature extraction, which carries over to the proposed sparse representation framework, is reduced data dimension and computational cost. For raw face images, the corresponding linear system  $y = Ax$  is very large. For instance, if the face images are given at the typical resolution,  $640 \times 480$  pixels, the dimension  $m$  is in the order of  $10^5$ . Although Algorithm 1 relies on scalable methods such as linear programming, directly applying it to such high-resolution images is still beyond the capability of regular computers.

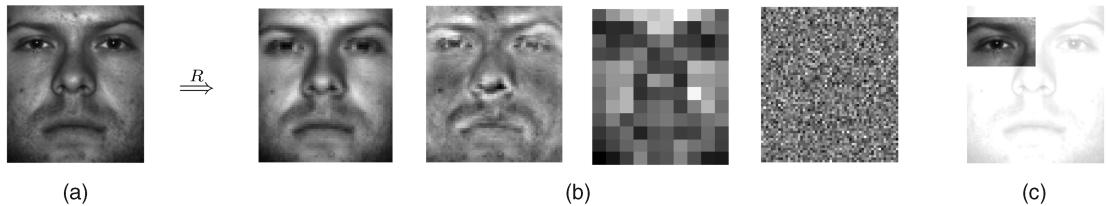
Since most feature transformations involve only linear operations (or approximately so), the projection from the image space to the feature space can be represented as a matrix  $R \in \mathbb{R}^{d \times m}$  with  $d \ll m$ . Applying  $R$  to both sides of (3) yields

$$\tilde{y} \doteq Ry = RAx_0 \in \mathbb{R}^d. \quad (16)$$

In practice, the dimension  $d$  of the feature space is typically chosen to be much smaller than  $n$ . In this case, the system of equations  $\tilde{y} = RAx \in \mathbb{R}^d$  is underdetermined in the unknown  $x \in \mathbb{R}^n$ . Nevertheless, as the desired solution  $x_0$  is sparse, we can hope to recover it by solving the following reduced  $\ell^1$ -minimization problem:

$$(\ell_r^1) : \hat{x}_1 = \arg \min \|x\|_1 \quad \text{subject to} \quad \|RAx - \tilde{y}\|_2 \leq \varepsilon, \quad (17)$$

for a given error tolerance  $\varepsilon > 0$ . Thus, in Algorithm 1, the matrix  $A$  of training images is now replaced by the matrix



**Fig. 6. Examples of feature extraction.** (a) Original face image. (b) 120D representations in terms of four different features (from left to right): Eigenfaces, Laplacianfaces, downsampled ( $12 \times 10$  pixel) image, and random projection. We will demonstrate that all these features contain almost the same information about the identity of the subject and give similarly good recognition performance. (c) The eye is a popular choice of feature for face recognition. In this case, the feature matrix  $R$  is simply a binary mask. (a) Original  $y$ . (b) 120D features  $\tilde{y} = Ry$ . (c) Eye feature  $y$ .

$RA \in \mathbb{R}^{d \times n}$  of  $d$ -dimensional features; the test image  $y$  is replaced by its features  $\tilde{y}$ .

For extant face recognition methods, empirical studies have shown that increasing the dimension  $d$  of the feature space generally improves the recognition rate, as long as the distribution of features  $RA_i$  does not become degenerate [42]. Degeneracy is not an issue for  $\ell^1$ -minimization, since it merely requires that  $\tilde{y}$  be in or near the range of  $RA_i$ —it does not depend on the covariance  $\Sigma_i = A_i^T R^T R A_i$  being nonsingular as in classical discriminant analysis. The stable version of  $\ell^1$ -minimization (10) or (17) is known in statistical literature as the Lasso [14].<sup>11</sup> It effectively regularizes highly underdetermined linear regression when the desired solution is sparse and has also been proven consistent in some noisy over-determined settings [12].

For our sparse representation approach to recognition, we would like to understand how the choice of the feature extraction  $R$  affects the ability of the  $\ell^1$ -minimization (17) to recover the correct sparse solution  $x_0$ . From the geometric interpretation of  $\ell^1$ -minimization given in Section 2.2.1, the answer to this depends on whether the associated new polytope  $P = RA(P_1)$  remains sufficiently neighborly. It is easy to show that the neighborliness of the polytope  $P = RA(P_1)$  increases with  $d$  [11], [15]. In Section 4, our experimental results will verify the ability of  $\ell^1$ -minimization, in particular, the stable version (17), to recover sparse representations for face recognition using a variety of features. This suggests that most data-dependent features popular in face recognition (e.g., eigenfaces and Laplacianfaces) may indeed give highly neighborly polytopes  $P$ .

Further analysis of high-dimensional polytope geometry has revealed a somewhat surprising phenomenon: if the solution  $x_0$  is sparse enough, then with overwhelming probability, it can be correctly recovered via  $\ell^1$ -minimization from any sufficiently large number  $d$  of linear measurements  $\tilde{y} = RAx_0$ . More precisely, if  $x_0$  has  $t \ll n$  nonzeros, then with overwhelming probability

$$d \geq 2t \log(n/d) \quad (18)$$

11. Classically, the Lasso solution is defined as the minimizer of  $\|y - Ax\|_2^2 + \lambda \|x\|_1$ . Here,  $\lambda$  can be viewed as inverse of the Lagrange multiplier associated with a constraint  $\|y - Ax\|_2^2 \leq \varepsilon$ . For every  $\lambda$ , there is an  $\varepsilon$  such that the two problems have the same solution. However,  $\varepsilon$  can be interpreted as a pixel noise level and fixed across various instances of the problem, whereas  $\lambda$  cannot. One should distinguish the Lasso optimization problem from the LARS algorithm, which provably solves some instances of Lasso with very sparse optimizers [35].

random linear measurements are sufficient for  $\ell^1$ -minimization (17) to recover the correct sparse solution  $x_0$  [44].<sup>12</sup> This surprising phenomenon has been dubbed the “blessing of dimensionality” [15], [46]. Random features can be viewed as a less-structured counterpart to classical face features such as Eigenfaces or Fisherfaces. Accordingly, we call the linear projection generated by a Gaussian random matrix *Randomfaces*.<sup>13</sup>

**Definition 2 (randomfaces).** Consider a transform matrix  $R \in \mathbb{R}^{d \times m}$  whose entries are independently sampled from a zero-mean normal distribution, and each row is normalized to unit length. The row vectors of  $R$  can be viewed as  $d$  random faces in  $\mathbb{R}^m$ .

One major advantage of Randomfaces is that they are extremely efficient to generate, as the transformation  $R$  is independent of the training data set. This advantage can be important for a face recognition system, where we may not be able to acquire a complete database of all subjects of interest to precompute data-dependent transformations such as Eigenfaces, or the subjects in the database may change over time. In such cases, there is no need for recomputing the random transformation  $R$ .

As long as the correct sparse solution  $x_0$  can be recovered, Algorithm 1 will always give the same classification result, regardless of the feature actually used. Thus, when the dimension of feature  $d$  exceeds the above bound (18), one should expect that the recognition performance of Algorithm 1 with different features quickly converges, and the choice of an “optimal” feature transformation is no longer critical: even random projections or downsampled images should perform as well as any other carefully engineered features. This will be corroborated by the experimental results in Section 4.

### 3.2 Robustness to Occlusion or Corruption

In many practical face recognition scenarios, the test image  $y$  could be partially corrupted or occluded. In this case, the above linear model (3) should be modified as

$$y = y_0 + e_0 = Ax_0 + e_0, \quad (19)$$

12. Strictly speaking, this threshold holds when random measurements are computed directly from  $x_0$ , i.e.,  $\tilde{y} = Rx_0$ . Nevertheless, our experiments roughly agree with the bound given by (18). The case where  $x_0$  is instead sparse in some overcomplete basis  $A$ , and we observe that random measurements  $\tilde{y} = RAx_0$  has also been studied in [45]. While conditions for correct recovery have been given, the bounds are not yet as sharp as (18) above.

13. Random projection has been previously studied as a general dimensionality-reduction method for numerous clustering problems [47], [48], [49], as well as for learning nonlinear manifolds [50], [51].

where  $e_0 \in \mathbb{R}^m$  is a vector of errors—a fraction,  $\rho$ , of its entries are nonzero. The nonzero entries of  $e_0$  model which pixels in  $y$  are corrupted or occluded. The locations of corruption can differ for different test images and are not known to the computer. The errors may have arbitrary magnitude and therefore cannot be ignored or treated with techniques designed for small noise such as the one given in Section 2.2.2.

A fundamental principle of coding theory [52] is that redundancy in the measurement is essential to detecting and correcting gross errors. Redundancy arises in object recognition because the number of image pixels is typically far greater than the number of subjects that have generated the images. In this case, even if a fraction of the pixels are completely corrupted by occlusion, recognition may still be possible based on the remaining pixels. On the other hand, feature extraction schemes discussed in the previous section would discard useful information that could help compensate for the occlusion. In this sense, no representation is more redundant, robust, or informative than the original images. Thus, when dealing with occlusion and corruption, we should always work with the highest possible resolution, performing downsampling or feature extraction only if the resolution of the original images is too high to process.

Of course, redundancy would be of no use without efficient computational tools for exploiting the information encoded in the redundant data. The difficulty in directly harnessing the redundancy in corrupted raw images has led researchers to instead focus on *spatial locality* as a guiding principle for robust recognition. Local features computed from only a small fraction of the image pixels are clearly less likely to be corrupted by occlusion than holistic features. In face recognition, methods such as ICA [53] and LNMF [54] exploit this observation by adaptively choosing filter bases that are locally concentrated. Local Binary Patterns [55] and Gabor wavelets [56] exhibit similar properties, since they are also computed from local image regions. A related approach partitions the image into fixed regions and computes features for each region [16], [57]. Notice, though, that projecting onto locally concentrated bases transforms the domain of the occlusion problem, rather than eliminating the occlusion. Errors on the original pixels become errors in the transformed domain and may even become less local. The role of feature extraction in achieving spatial locality is therefore questionable, since *no bases or features are more spatially localized than the original image pixels themselves*. In fact, the most popular approach to robustifying feature-based methods is based on randomly sampling individual pixels [28], sometimes in conjunction with statistical techniques such as multivariate trimming [29].

Now, let us show how the proposed sparse representation classification framework can be extended to deal with occlusion. Let us assume that the corrupted pixels are a relatively small portion of the image. The error vector  $e_0$ , like the vector  $x_0$ , then has sparse<sup>14</sup> nonzero entries. Since  $y_0 = Ax_0$ , we can rewrite (19) as

14. Here, “sparse” does not mean “very few.” In fact, as our experiments will demonstrate, the portion of corrupted entries can be rather significant. Depending on the type of corruption, our method can handle up to  $\rho = 40$  percent or  $\rho = 70$  percent corrupted pixels.

$$\mathbf{y} = [\mathbf{A}, \mathbf{I}] \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{e}_0 \end{bmatrix} = \mathbf{B}\mathbf{w}_0. \quad (20)$$

Here,  $\mathbf{B} = [\mathbf{A}, \mathbf{I}] \in \mathbb{R}^{m \times (n+m)}$ , so the system  $\mathbf{y} = \mathbf{B}\mathbf{w}$  is always underdetermined and does not have a unique solution for  $\mathbf{w}$ . However, from the above discussion about the sparsity of  $\mathbf{x}_0$  and  $\mathbf{e}_0$ , the correct generating  $\mathbf{w}_0 = [\mathbf{x}_0, \mathbf{e}_0]$  has at most  $n_i + \rho m$  nonzeros. We might therefore hope to recover  $\mathbf{w}_0$  as the sparsest solution to the system  $\mathbf{y} = \mathbf{B}\mathbf{w}$ . In fact, if the matrix  $\mathbf{B}$  is in general position, then as long as  $\mathbf{y} = \mathbf{B}\tilde{\mathbf{w}}$  for some  $\tilde{\mathbf{w}}$  with less than  $m/2$  nonzeros,  $\tilde{\mathbf{w}}$  is the unique sparsest solution. Thus, if the occlusion  $\mathbf{e}$  covers less than  $\frac{m-n_i}{2}$  pixels,  $\approx 50$  percent of the image, the sparsest solution  $\tilde{\mathbf{w}}$  to  $\mathbf{y} = \mathbf{B}\mathbf{w}$  is the true generator,  $\mathbf{w}_0 = [\mathbf{x}_0, \mathbf{e}_0]$ .

More generally, one can assume that the corrupting error  $\mathbf{e}_0$  has a sparse representation with respect to some basis  $\mathbf{A}_e \in \mathbb{R}^{m \times n_e}$ . That is,  $\mathbf{e}_0 = \mathbf{A}_e \mathbf{u}_0$  for some sparse vector  $\mathbf{u}_0 \in \mathbb{R}^m$ . Here, we have chosen the special case  $\mathbf{A}_e = \mathbf{I} \in \mathbb{R}^{m \times m}$  as  $\mathbf{e}_0$  is assumed to be sparse with respect to the natural pixel coordinates. If the error  $\mathbf{e}_0$  is instead more sparse with respect to another basis, e.g., Fourier or Haar, we can simply redefine the matrix  $\mathbf{B}$  by appending  $\mathbf{A}_e$  (instead of the identity  $\mathbf{I}$ ) to  $\mathbf{A}$  and instead seek the sparsest solution  $\mathbf{w}_0$  to the equation:

$$\mathbf{y} = \mathbf{B}\mathbf{w} \quad \text{with} \quad \mathbf{B} = [\mathbf{A}, \mathbf{A}_e] \in \mathbb{R}^{m \times (n+n_e)}. \quad (21)$$

In this way, the same formulation can handle more general classes of (sparse) corruption.

As before, we attempt to recover the sparsest solution  $\mathbf{w}_0$  from solving the following *extended  $\ell^1$ -minimization* problem:

$$(\ell_e^1) : \quad \hat{\mathbf{w}}_1 = \arg \min \|\mathbf{w}\|_1 \quad \text{subject to} \quad \mathbf{B}\mathbf{w} = \mathbf{y}. \quad (22)$$

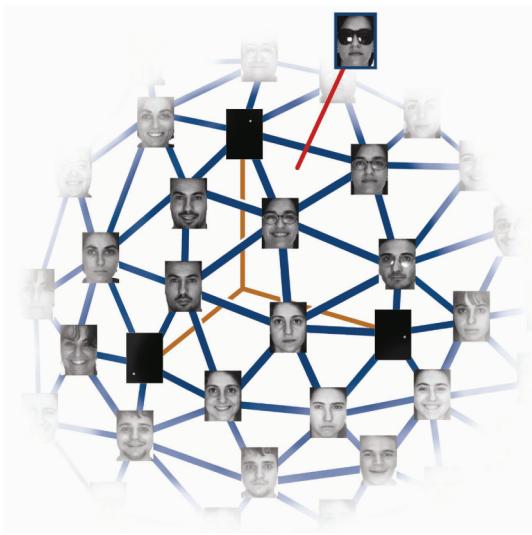
That is, in Algorithm 1, we now replace the image matrix  $\mathbf{A}$  with the extended matrix  $\mathbf{B} = [\mathbf{A}, \mathbf{I}]$  and  $\mathbf{x}$  with  $\mathbf{w} = [\mathbf{x}, \mathbf{e}]$ .

Clearly, whether the sparse solution  $\mathbf{w}_0$  can be recovered from the above  $\ell^1$ -minimization depends on the neighborliness of the new polytope  $P = \mathbf{B}(P_1) = [\mathbf{A}, \mathbf{I}](P_1)$ . This polytope contains vertices from both the training images  $\mathbf{A}$  and the identity matrix  $\mathbf{I}$ , as illustrated in Fig. 7. The bounds given in (8) imply that if  $\mathbf{y}$  is an image of subject  $i$ , the  $\ell^1$ -minimization (22) cannot guarantee to correctly recover  $\mathbf{w}_0 = [\mathbf{x}_0, \mathbf{e}_0]$  if

$$n_i + |\text{support}(\mathbf{e}_0)| > d/3.$$

Generally,  $d \gg n_i$ , so, (8) implies that the largest fraction of occlusion under which we can hope to still achieve perfect reconstruction is 33 percent. This bound is corroborated by our experimental results, see Fig. 12.

To know exactly how much occlusion can be tolerated, we need more accurate information about the neighborliness of the polytope  $P$  than a loose upper bound given by (8). For instance, we would like to know for a given set of training images, what is the largest amount of (worst possible) occlusion it can handle. While the best known algorithms for exactly computing the neighborliness of a polytope are combinatorial in nature, tighter upper bounds can be obtained by restricting the search for intersections between the nullspace of  $\mathbf{B}$  and the  $\ell^1$ -ball to a random subset of the  $t$ -faces of the  $\ell^1$ -ball (see [37] for details). We



**Fig. 7. Face recognition with occlusion.** The columns of  $\pm B = \pm[A, I]$  span a high-dimensional polytope  $P = B(P_1)$  in  $\mathbb{R}^n$ . Each vertex of this polytope is either a training image or an image with just a single pixel illuminated (corresponding to the identity submatrix  $I$ ). Given a test image, solving the  $\ell^1$ -minimization problem essentially locates which facet of the polytope the test image falls on. The  $\ell^1$ -minimization finds the facet with the fewest possible vertices. Only vertices of that facet contribute to the representation; all other vertices have no contribution.

will use this technique to estimate the neighborliness of all the training data sets considered in our experiments.

Empirically, we found that the stable version (10) is only necessary when we do not consider occlusion or corruption  $e_0$  in the model (such as the case with feature extraction discussed in the previous section). When we explicitly account for gross errors by using  $B = [A, I]$  the extended  $\ell^1$ -minimization (22) with the exact constraint  $Bw = y$  is already stable under moderate noise.

Once the sparse solution  $\hat{w}_1 = [\hat{x}_1, \hat{e}_1]$  is computed, setting  $y_r = y - \hat{e}_1$  recovers a clean image of the subject with occlusion or corruption compensated for. To identify the subject, we slightly modify the residual  $r_i(y)$  in Algorithm 1, computing it against the recovered image  $y_r$ :

$$r_i(y) = \|y_r - A\delta_i(\hat{x}_1)\|_2 = \|y - \hat{e}_1 - A\delta_i(\hat{x}_1)\|_2. \quad (23)$$

## 4 EXPERIMENTAL VERIFICATION

In this section, we present experiments on publicly available databases for face recognition, which serve both to demonstrate the efficacy of the proposed classification algorithm and to validate the claims of the previous sections. We will first examine the role of feature extraction within our framework, comparing performance across various feature spaces and feature dimensions, and comparing to several popular classifiers. We will then demonstrate the robustness of the proposed algorithm to corruption and occlusion. Finally, we demonstrate (using ROC curves) the effectiveness of sparsity as a means of validating test images and examine how to choose training sets to maximize robustness to occlusion.

### 4.1 Feature Extraction and Classification Methods

We test our SRC algorithm using several conventional holistic face features, namely, Eigenfaces, Laplacianfaces, and Fisherfaces, and compare their performance with two

unconventional features: randomfaces and downsampled images. We compare our algorithm with three classical algorithms, namely, NN, and NS, discussed in the previous section, as well as linear SVM.<sup>15</sup> In this section, we use the stable version of SRC in various lower dimensional feature spaces, solving the reduced optimization problem (17) with the error tolerance  $\epsilon = 0.05$ . The Matlab implementation of the reduced (feature space) version of Algorithm 1 takes only a few seconds per test image on a typical 3-GHz PC.

#### 4.1.1 Extended Yale B Database

The Extended Yale B database consists of 2,414 frontal-face images of 38 individuals [58]. The cropped and normalized  $192 \times 168$  face images were captured under various laboratory-controlled lighting conditions [59]. For each subject, we randomly select half of the images for training (i.e., about 32 images per subject) and the other half for testing. Randomly choosing the training set ensures that our results and conclusions will not depend on any special choice of the training data.

We compute the recognition rates with the feature space dimensions 30, 56, 120, and 504. Those numbers correspond to downsampling ratios of 1/32, 1/24, 1/16, and 1/8, respectively.<sup>16</sup> Notice that Fisherfaces are different from the other features because the maximal number of valid Fisherfaces is one less than the number of classes  $k$  [24], 38 in this case. As a result, the recognition result for Fisherfaces is only available at dimension 30 in our experiment.

The subspace dimension for the NS algorithm is 9, which has been mostly agreed upon in the literature for processing facial images with only illumination change.<sup>17</sup> Fig. 8 shows the recognition performance for the various features, in conjunction with four different classifiers: SRC, NN, NS, and SVM.

SRC achieves recognition rates between 92.1 percent and 95.6 percent for all 120D feature spaces and a maximum rate of 98.1 percent with 504D randomfaces.<sup>18</sup> The maximum recognition rates for NN, NS, and SVM are 90.7 percent, 94.1 percent, and 97.7 percent, respectively. Tables with all the recognition rates are available in the supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>. The recognition rates shown in Fig. 8 are consistent with those that have been reported in the literature, although some reported on different databases or with different training subsets. For example, He et al. [25] reported the best recognition rate of 75 percent using Eigenfaces at 33D, and 89 percent using Laplacianfaces at

15. Due to the subspace structure of face images, linear SVM is already appropriate for separating features from different faces. The use of a linear kernel (as opposed to more complicated nonlinear transformations) also makes it possible to directly compare between different algorithms working in the same feature space. Nevertheless, better performance might be achieved by using nonlinear kernels in addition to feature transformations.

16. We cut off the dimension at 504 as the computation of Eigenfaces and Laplacianfaces reaches the memory limit of Matlab. Although our algorithm persists to work far beyond on the same computer, 504 is already sufficient to reach all our conclusions.

17. Subspace dimensions significantly greater or less than 9 eventually led to a decrease in performance.

18. We also experimented with replacing the constrained  $\ell^1$ -minimization in the SRC algorithm with the Lasso. For appropriate choice of regularization  $\lambda$ , the results are similar. For example, with downsampled faces as features and  $\lambda = 1,000$ , the recognition rates are 73.7 percent, 86.2 percent, 91.9 percent, 97.5 percent, at dimensions 30, 56, 120, and 504 (within 1 percent of the results in Fig. 8).

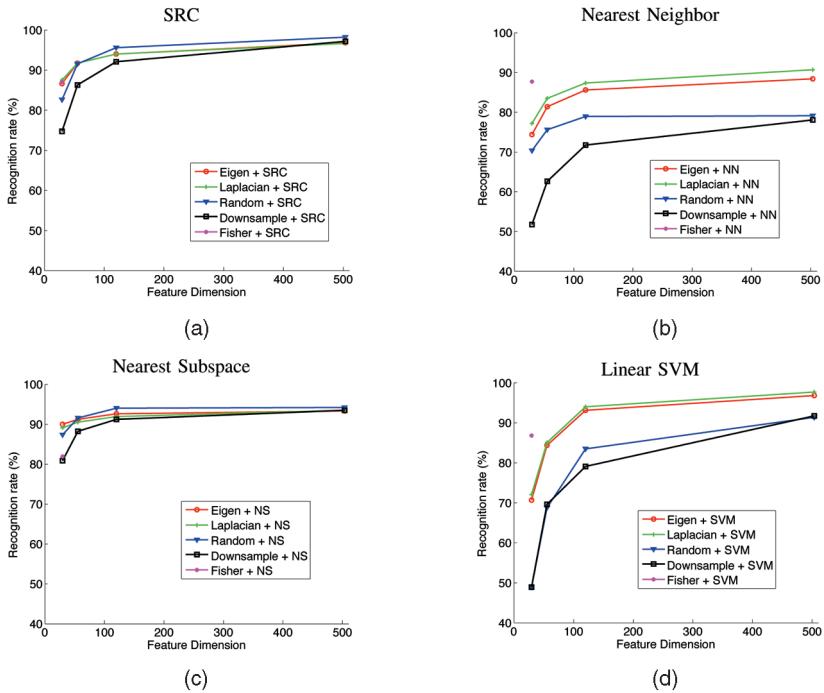


Fig. 8. Recognition rates on Extended Yale B database, for various feature transformations and classifiers. (a) SRC (our approach). (b) NN. (c) NS. (d) SVM (linear kernel).

28D on the Yale face database, both using NN. In [32], Lee et al. reported 95.4 percent accuracy using the NS method on the Yale B database.

#### 4.1.2 AR Database

The AR database consists of over 4,000 frontal images for 126 individuals. For each individual, 26 pictures were taken in two separate sessions [60]. These images include more facial variations, including illumination change, expressions, and facial disguises comparing to the Extended Yale B database. In the experiment, we chose a subset of the data set consisting of 50 male subjects and 50 female subjects. For each subject, 14 images with only illumination change and expressions were selected: the seven images from Session 1 for training, and the other seven from Session 2 for testing. The images are cropped with dimension  $165 \times 120$  and converted to gray scale. We selected four feature space dimensions: 30, 54, 130, and 540, which correspond to the downsample ratios 1/24, 1/18, 1/12, and 1/6, respectively. Because the number of subjects is 100, results for Fisherfaces are only given at dimension 30 and 54.

This database is substantially more challenging than the Yale database, since the number of subjects is now 100, but the training images is reduced to seven per subject: four neutral faces with different lighting conditions and three faces with different expressions. For NS, since the number of training images per subject is seven, any estimate of the face subspace cannot have dimension higher than 7. We chose to keep all seven dimensions for NS in this case.

Fig. 9 shows the recognition rates for this experiment. With 540D features, SRC achieves a recognition rate between 92.0 percent and 94.7 percent. On the other hand, the best rates achieved by NN and NS are 89.7 percent and 90.3 percent, respectively. SVM slightly outperforms SRC on this data set, achieving a maximum recognition rate of 95.7 percent. However, the performance of SVM varies more

with the choice of feature space—the recognition rate using random features is just 88.8 percent. The supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>, contains a table of detailed numerical results.

Based on the results on the Extended Yale B database and the AR database, we draw the following conclusions:

- For both the Yale database and AR database, the best performances of SRC and SVM consistently exceed the best performances of the two classical methods NN and NS at each individual feature dimension. More specifically, the best recognition rate for SRC on the Yale database is 98.1 percent, compared to 97.7 percent for SVM, 94.0 percent for NS, and 90.7 percent for NN; the best rate for SRC on the AR database is 94.7 percent, compared to 95.7 percent for SVM, 90.3 percent for NS, and 89.7 percent for NN.
- The performances of the other three classifiers depends strongly on a good choice of “optimal” features—Fisherfaces for lower feature space dimension and Laplacianfaces for higher feature space dimension. With NN and SVM, the performance of the various features does not converge as the dimension of the feature space increases.
- The results corroborate the theory of compressed sensing: (18) suggests that  $d \approx 128$  random linear measurements should suffice for sparse recovery in the Yale database, while  $d \approx 88$  random linear measurements should suffice for sparse recovery in the AR database [44]. Beyond these dimensions, the performances of various features in conjunction with  $\ell^1$ -minimization converge, with conventional and unconventional features (e.g., Randomfaces and downsampled images) performing similarly. When

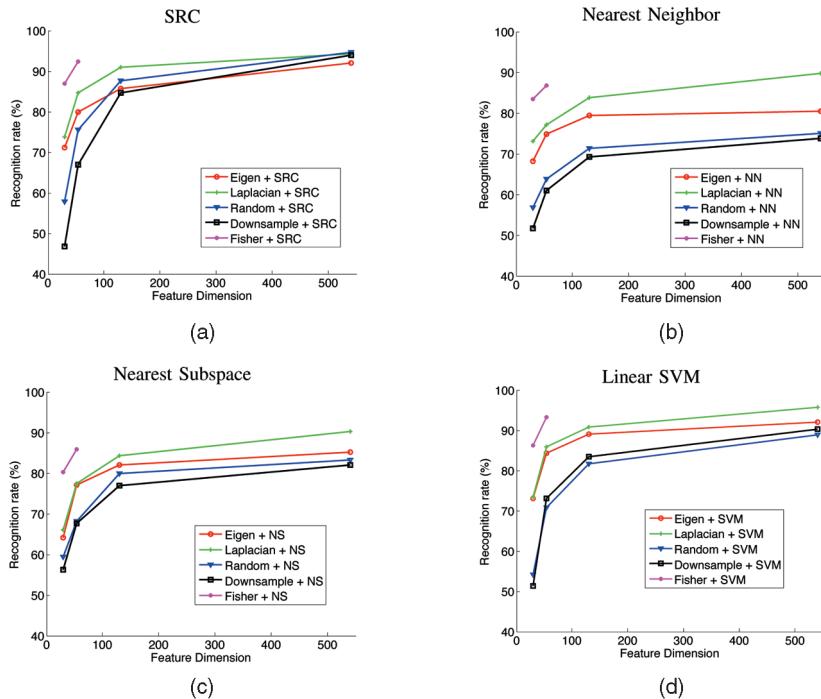


Fig. 9. Recognition rates on AR database, for various feature transformations and classifiers. (a) SRC (our approach). (b) NN. (c) NS. (d) SVM (linear kernel).

the feature dimension is large, a single random projection performs the best (98.1 percent recognition rate on Yale, 94.7 percent on AR).

## 4.2 Partial Face Features

here have been extensive studies in both the human and computer vision literature about the effectiveness of partial features in recovering the identity of a human face, e.g., see [21] and [41]. As a second set of experiments, we test our algorithm on the following three partial facial features: nose, right eye, and mouth and chin. We use the Extended Yale B database for the experiment, with the same training and test sets, as in Section 4.1.1. See Fig. 10 for a typical example of the extracted features.

For each of the three features, the dimension  $d$  is larger than the number of training samples ( $n = 1,207$ ), and the linear system (16) to be solved becomes overdetermined. Nevertheless, sparse approximate solutions  $x$  can still be

obtained by solving the  $\varepsilon$ -relaxed  $\ell^1$ -minimization problem (17) (here, again,  $\varepsilon = 0.05$ ). The results in Fig. 10 right again show that the proposed SRC algorithm achieves better recognition rates than NN, NS, and SVM. These experiments also show the scalability of the proposed algorithm in working with more than  $10^4$ -dimensional features.

## 4.3 Recognition Despite Random Pixel Corruption

For this experiment, we test the robust version of SRC, which solves the extended  $\ell^1$ -minimization problem (22) using the Extended Yale B Face Database. We choose Subsets 1 and 2 (717 images, normal-to-moderate lighting conditions) for training and Subset 3 (453 images, more extreme lighting conditions) for testing. Without occlusion, this is a relatively easy recognition problem. This choice is deliberate, in order to isolate the effect of occlusion. The images are resized to  $96 \times 84$  pixels,<sup>19</sup> so in this case,  $B = [A, I]$  is an  $8,064 \times 8,761$  matrix. For this data set, we have estimated that the polytope  $P = \text{conv}(\pm B)$  is approximately 1,185 neighborly (using the method given in [37]), suggesting that perfect reconstruction can be achieved up to 13.3 percent (worst possible) occlusion.

We corrupt a percentage of randomly chosen pixels from each of the test images, replacing their values with independent and identically distributed samples from a uniform distribution.<sup>20</sup> The corrupted pixels are randomly chosen for each test image, and the locations are unknown to the algorithm. We vary the percentage of corrupted pixels from 0 percent to 90 percent. Figs. 11a, 11b, 11c, and 11d shows several example test images. To the human eye, beyond 50 percent corruption, the corrupted images (Fig. 11a second and third rows) are

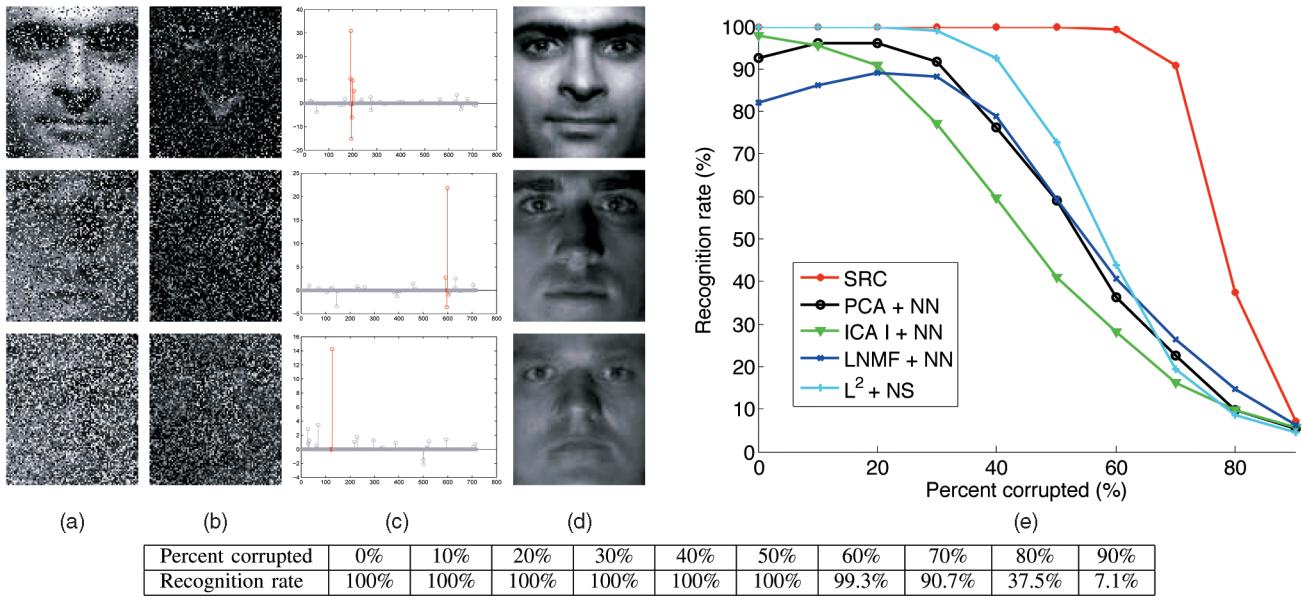


Features	Nose	Right Eye	Mouth & Chin
Dimension ( $d$ )	4,270	5,040	12,936
SRC	<b>87.3%</b>	<b>93.7%</b>	<b>98.3%</b>
NN	49.2%	68.8%	72.7%
NS	83.7%	78.6%	94.4%
SVM	70.8%	85.8%	95.3%

Fig. 10. Recognition with partial face features. (a) Example features. (b) Recognition rates of SRC, NN, NS, and SVM on the Extended Yale B database.

19. The only reason for resizing the images is to be able to run all the experiments within the memory size of Matlab on a typical PC. The algorithm relies on linear programming and is scalable in the image size.

20. Uniform over  $[0, y_{\max}]$ , where  $y_{\max}$  is the largest possible pixel value.



**Fig. 11. Recognition under random corruption.** (a) Test images  $y$  from Extended Yale B, with random corruption. Top row: 30 percent of pixels are corrupted. Middle row: 50 percent corrupted. Bottom row: 70 percent corrupted. (b) Estimated errors  $\hat{e}_1$ . (c) Estimated sparse coefficients  $\hat{x}_1$ . (d) Reconstructed images  $y_r$ . SRC correctly identifies all three corrupted face images. (e) The recognition rate across the entire range of corruption for various algorithms. SRC (red curve) significantly outperforms others, performing almost perfectly up to 60 percent random corruption (see table below).

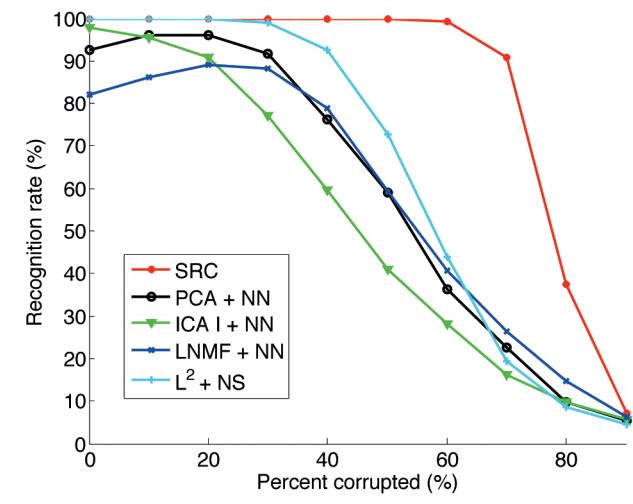
barely recognizable as face images; determining their identity seems out of the question. Nevertheless, even in this extreme circumstance, SRC correctly recovers the identity of the subjects.

We quantitatively compare our method to four popular techniques for face recognition in the vision literature. The Principal Component Analysis (PCA) approach in [23] is not robust to occlusion. There are many variations to make PCA robust to corruption or incomplete data, and some have been applied to robust face recognition, e.g., [29]. We will later discuss their performance against ours on more realistic conditions. Here, we use the basic PCA to provide a standard baseline for comparison.<sup>21</sup> The remaining three techniques are designed to be more robust to occlusion. Independent Component Analysis (ICA) architecture I [53] attempts to express the training set as a linear combination of statistically independent basis images. Local Nonnegative Matrix Factorization (LNMF) [54] approximates the training set as an additive combination of basis images, computed with a bias toward sparse bases.<sup>22</sup> Finally, to demonstrate that the improved robustness is really due to the use of the  $\ell^1$ -norm, we compare to a least-squares technique that first projects the test image onto the subspace spanned by all face images and then performs NS.

Fig. 11e plots the recognition performance of SRC and its five competitors, as a function of the level of corruption. We see that the algorithm dramatically outperforms others. From 0 percent upto 50 percent occlusion, SRC correctly classifies all subjects. At 50 percent corruption, none of the others achieves higher than 73 percent recognition rate, while the proposed algorithm achieves 100 percent. Even at 70 percent occlusion, the recognition rate is still 90.7 percent.

21. Following [58], we normalize the image pixels to have zero mean and unit variance before applying PCA.

22. For PCA, ICA, and LNMF, the number of basis components is chosen to give the optimal test performance over the range  $\{100, 200, 300, 400, 500, 600\}$ .

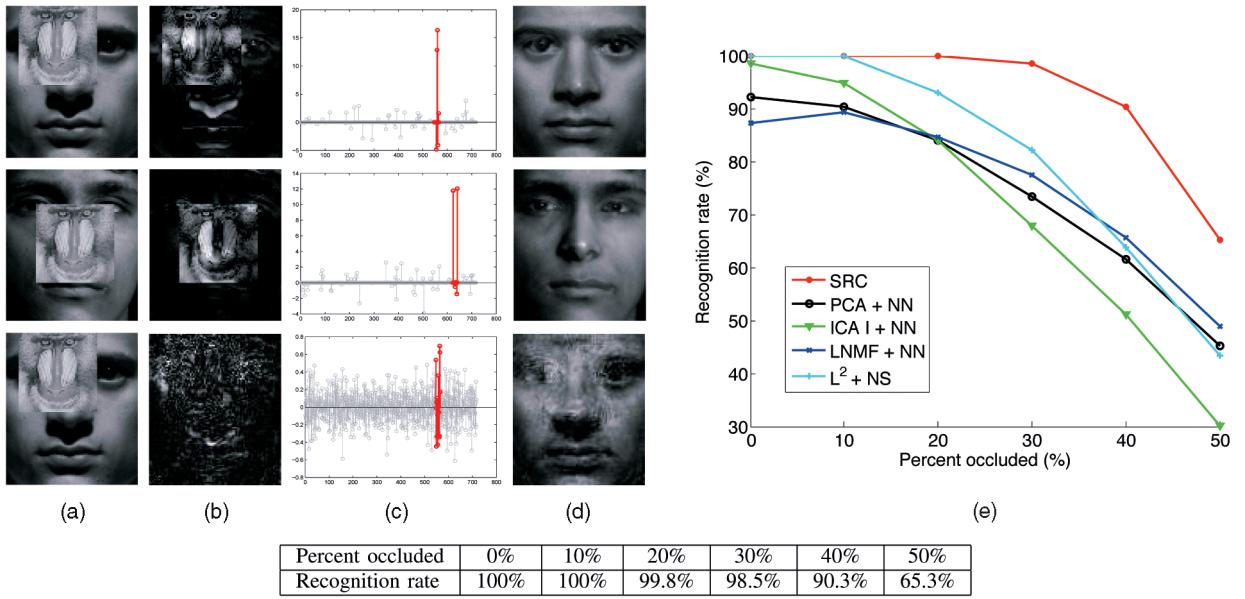


This greatly surpasses the theoretical bound of the worst-case corruption (13.3 percent) that the algorithm is ensured to tolerate. Clearly, the worst-case analysis is too conservative for random corruption.

#### 4.4 Recognition Despite Random Block Occlusion

We next simulate various levels of contiguous occlusion, from 0 percent to 50 percent, by replacing a randomly located square block of each test image with an unrelated image, as in Fig. 12a. Again, the location of occlusion is randomly chosen for each image and is unknown to the computer. Methods that select fixed facial features or blocks of the image (e.g., [16] and [57]) are less likely to succeed here due to the unpredictable location of the occlusion. The top two rows in Figs. 12a, 12b, 12c, and 12d shows the two representative results of Algorithm 1 with 30 percent occlusion. Fig. 12a is the occluded image. In the second row, the entire center of the face is occluded; this is a difficult recognition task even for humans. Fig. 12b shows the magnitude of the estimated error  $\hat{e}_1$ . Notice that  $\hat{e}_1$  compensates not only for occlusion due to the baboon but also for the violation of the linear subspace model caused by the shadow under the nose. Fig. 12c plots the estimated coefficient vector  $\hat{x}_1$ . The red entries are coefficients corresponding to test image's true class. In both examples, the estimated coefficients are indeed sparse and have large magnitude only for training images of the same person. In both cases, the SRC algorithm correctly classifies the occluded image. For this data set, our Matlab implementation requires 90 seconds per test image on a PowerMac G5.

The graph in Fig. 12e shows the recognition rates of all six algorithms. SRC again significantly outperforms the other five methods for all levels of occlusion. Upto 30 percent occlusion, Algorithm 1 performs almost perfectly, correctly identifying over 98 percent of test subjects. Even at 40 percent occlusion, only 9.7 percent of subjects are misclassified. Compared to the random pixel corruption, contiguous



**Fig. 12. Recognition under varying level of contiguous occlusion.** Left, top two rows: (a) 30 percent occluded test face images  $y$  from Extended Yale B. (b) Estimated sparse errors,  $\hat{e}_1$ . (c) Estimated sparse coefficients,  $\hat{x}_1$ , red (darker) entries correspond to training images of the same person. (d) Reconstructed images,  $y_r$ . SRC correctly identifies both occluded faces. For comparison, the bottom row shows the same test case, with the result given by least squares (overdetermined  $\ell^2$ -minimization). (e) The recognition rate across the entire range of corruption for various algorithms. SRC (red curve) significantly outperforms others, performing almost perfectly up to 30 percent contiguous occlusion (see table below).

occlusion is certainly a worse type of errors for the algorithm. Notice, though, that the algorithm does not assume any knowledge about the nature of corruption or occlusion. In Section 4.6, we will see how prior knowledge that the occlusion is contiguous can be used to customize the algorithm and greatly enhance the recognition performance.

This result has interesting implications for the debate over the use of holistic versus local features in face recognition [22]. It has been suggested that both ICA I and LNMF are robust to occlusion: since their bases are locally concentrated, occlusion corrupts only a fraction of the coefficients. By contrast, if one uses  $\ell^2$ -minimization (orthogonal projection) to express an occluded image in terms of a holistic basis such as the training images themselves, all of the coefficients may be corrupted (as in Fig. 12 third row). The implication here is that the problem is *not* the choice of representing the test image in terms of a holistic or local basis, but rather *how the representation is computed*. Properly harnessing redundancy and sparsity is the key to error correction and robustness. Extracting local or disjoint features can only reduce redundancy, resulting in inferior robustness.

#### 4.5 Recognition Despite Disguise

We test SRC's ability to cope with real possibly malicious occlusions using a subset of the ARFace Database. The chosen subset consists of 1,399 images (14 each, except for a corrupted image w-027-14.bmp) of 100 subjects, 50 male and 50 female. For training, we use 799 images (about 8 per subject) of unoccluded frontal views with varying facial expression, giving a matrix  $B$  of size  $4,980 \times 5,779$ . We estimate  $P = \text{conv}(\pm B)$  is approximately 577 neighborly, indicating that perfect reconstruction is possible up to 11.6 percent occlusion. Our Matlab implementation requires about 75 seconds per test image on a PowerMac G5.

We consider two separate test sets of 200 images. The first test set contains images of the subjects wearing sunglasses, which occlude roughly 20 percent of the image.

Fig. 1a shows a successful example from this test set. Notice that  $\hat{e}_1$  compensates for small misalignment of the image edges, as well as occlusion due to sunglasses. The second test set considered contains images of the subjects wearing a scarf, which occludes roughly 40 percent of the image. Since the occlusion level is more than three times the maximum worst case occlusion given by the neighborliness of  $\text{conv}(\pm B)$ , our approach is unlikely to succeed in this domain. Fig. 13a shows one such failure. Notice that the largest coefficient corresponds to an image of a bearded man whose mouth region resembles the scarf.

The table in Fig. 13 left compares SRC to the other five algorithms described in the previous section. On faces occluded by sunglasses, SRC achieves a recognition rate of 87 percent, more than 17 percent better than the nearest competitor. For occlusion by scarves, its recognition rate is 59.5 percent, more than double its nearest competitor but still quite poor. This confirms that although the algorithm is provably robust to arbitrary occlusions upto the breakdown point determined by the neighborliness of the training set, beyond that point, it is sensitive to occlusions that resemble regions of a training image from a different individual. Because the amount of occlusion exceeds this breakdown point, additional assumptions, such as the disguise is likely to be contiguous, are needed to achieve higher recognition performance.

#### 4.6 Improving Recognition by Block Partitioning

Thus far, we have not exploited the fact that in many real recognition scenarios, the occlusion falls on some patch of image pixels which is a priori unknown but is known to be connected. A somewhat traditional approach (explored in [57] among others) to exploiting this information in face recognition is to partition the image into blocks and process each block independently. The results for individual blocks are then aggregated, for example, by voting, while discarding blocks believed to be occluded (using, say, the outlier

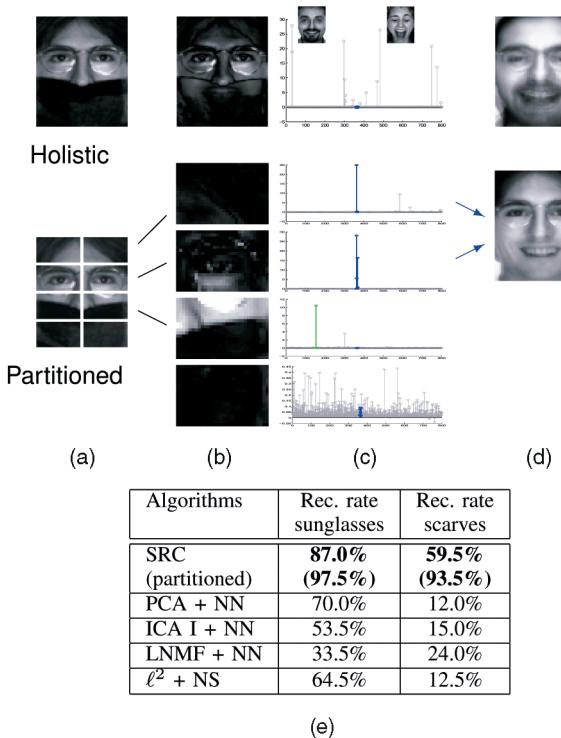


Fig. 13. (a)-(d) **Partition scheme to tackle contiguous disguise.** The top row visualizes an example for which SRC failed with the whole image (holistic). The two largest coefficients correspond to a bearded man and a screaming woman, two images whose mouth region resembles the occluding scarf. If the occlusion is known to be contiguous, one can partition the image into multiple smaller blocks, apply the SRC algorithm to each of the blocks and then aggregate the results by voting. The second row visualizes how this partition-based scheme works on the same test image but leads to a correct identification. (a) The test image, occluded by scarf. (b) Estimated sparse error  $\hat{\epsilon}_1$ . (c) Estimated sparse coefficients  $\hat{x}_1$ . (d) Reconstructed image. (e) **Table of recognition rates on the AR database.** The table shows the performance of all the algorithms for both types of occlusion. SRC, its holistic version (right top) and partitioned version (right bottom), achieves the highest recognition rate.

rejection rule introduced in Section 2.4). The major difficulty with this approach is that the occlusion cannot be expected to respect any fixed partition of the image; while only a few blocks are assumed to be completely occluded, some or all of the remaining blocks may be partially occluded. Thus, in such a scheme, there is still a need for robust techniques *within each block*.

We partition each of the training images into  $L$  blocks of size  $a \times b$ , producing a set of matrices  $A^{(1)}, \dots, A^{(L)} \in \mathbb{R}^{p \times n}$ , where  $p = ab$ . We similarly partition the test image  $y$  into  $y^{(1)}, \dots, y^{(L)} \in \mathbb{R}^p$ . We write the  $l$ th block of the test image as a sparse linear combination  $A^{(l)}x^{(l)}$  of  $l$ th blocks of the training images, plus a sparse error  $e^{(l)} \in \mathbb{R}^p : y^{(l)} = A^{(l)}x^{(l)} + e^{(l)}$ . We can recover can again recover a sparse  $w^{(l)} = [x^{(l)} \ e^{(l)}] \in \mathbb{R}^{n+p}$  by  $\ell^1$  minimization:

$$\hat{w}_1^{(l)} \doteq \arg \min_{w \in \mathbb{R}^{n+p}} \|w\|_1 \quad \text{subject to} \quad \begin{bmatrix} A^{(l)} & I \end{bmatrix} w = y^{(l)}. \quad (24)$$

We apply the classifier from Algorithm 1 within each block<sup>23</sup> and then aggregate the results by voting. Fig. 13 illustrates this scheme.

23. Occluded blocks can also be rejected via (15). We find that this does not significantly increase the recognition rate.

We verify the efficacy of this scheme on the AR database for faces disguised with sunglasses or scarves. We partition the images into eight ( $4 \times 2$ ) blocks of size  $20 \times 30$  pixels. Partitioning increases the recognition rate on scarves from 59.5 percent to 93.5 percent and also improves the recognition rate on sunglasses from 87.0 percent to 97.5 percent. This performance exceeds the best known results on the AR data set [29] to date. That work obtains 84 percent on the sunglasses and 93 percent on the scarfs, on only 50 subjects, using more sophisticated random sampling techniques. Also noteworthy is [16], which aims to recognize occluded faces from only a single training sample per subject. On the AR database, that method achieves a lower combined recognition rate of 80 percent.<sup>24</sup>

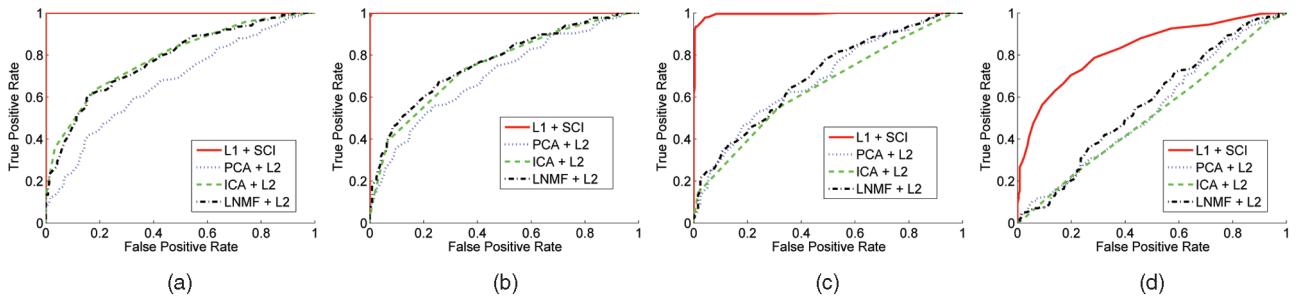
#### 4.7 Rejecting Invalid Test Images

We next demonstrate the relevance of sparsity for rejecting invalid test images, with or without occlusion. We test the outlier rejection rule (15) based on the Sparsity Concentration Index (14) on the Extended Yale B database, using Subsets 1 and 2 for training and Subset 3 for testing as before. We again simulate varying levels of occlusion (10 percent, 30 percent, and 50 percent) by replacing a randomly chosen block of each test image with an unrelated image. However, in this experiment, we include only half of the subjects in the training set. Thus, half of the subjects in the testing set are new to the algorithm. We test the system's ability to determine whether a given test subject is in the training database or not by sweeping the threshold  $\tau$  through a range of values in  $[0, 1]$ , generating the receiver operator characteristic (ROC) curves in Fig. 14. For comparison, we also considered outlier rejection by thresholding the euclidean distance between (features of) the test image and (features of) the nearest training images within the PCA, ICA, and LNMF feature spaces. These curves are also displayed in Fig. 14. Notice that the simple rejection rule (15) performs nearly perfectly at 10 percent and 30 percent occlusion. At 50 percent occlusion, it still significantly outperforms the other three algorithms and is the only one of the four algorithms that performs significantly better than chance. The supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>, contains more validation results on the AR database using Eigenfaces, again demonstrating significant improvement in the ROC.

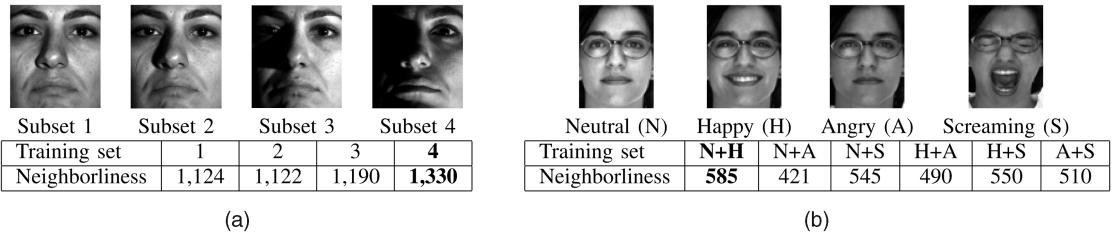
#### 4.8 Designing the Training Set for Robustness

An important consideration in designing recognition systems is selecting the number of training images, as well as the conditions (lighting, expression, viewpoint, etc.) under which they are to be taken. The training images should be extensive enough to span the conditions that might occur in the test set: they should be "sufficient" from a pattern recognition standpoint. For instance, Lee et al. [59] shows how to choose the fewest representative images to well approximate the illumination cone of each face. The notion of neighborliness discussed in Section 2 provides a different quantitative measure for how "robust" the training set is: the amount of worst case occlusion the algorithm can tolerate is directly determined by how neighborly the associated polytope is. The worst case is relevant in visual recognition,

24. From our own implementation and experiments, we find their method does not generalize well to more extreme illuminations.



**Fig. 14. ROC curves for outlier rejection.** Vertical axis: true positive rate. Horizontal axis: false positive rate. The solid red curve is generated by SRC with outliers rejected based on (15). The SCI-based validation and SRC classification together perform almost perfectly for up to 30 percent occlusion. (a) No occlusion. (b) Ten percent occlusion. (c) Thirty percent. (d) Fifty percent.



**Fig. 15. Robust training set design.** (a) Varying illumination. Top left: four subsets of Extended Yale B, containing increasingly extreme lighting conditions. Bottom left: estimated neighborliness of the polytope  $\text{conv}(\pm B)$  for each subset. (b) Varying expression. Top right: four facial expressions in the AR database. Bottom right: estimated neighborliness of  $\text{conv}(\pm B)$  when taking the training set from different pairs of expressions.

since the occluding object could potentially be quite similar to one of the other training classes. However, if the occlusion is random and uncorrelated with the training images, as in Section 4.3, the average behavior may also be of interest.

In fact, these two concerns, sufficiency and robustness, are complementary. Fig. 15a shows the estimated neighborliness for the four subsets of the Extended Yale B database. Notice that the highest neighborliness,  $\approx 1,330$ , is achieved with Subset 4, the most extreme lighting conditions. Fig. 15b shows the breakdown point for subsets of the AR database with different facial expressions. The data set contains four facial expressions, Neutral, Happy, Angry, and Scream, pictured in Fig. 15b. We generate training sets from all pairs of expressions and compute the neighborliness of each of the corresponding polytopes. The most robust training sets are achieved by the Neutral+Happy and Happy+Scream combinations, while the least robustness comes from Neutral+Angry. Notice that the Neutral and Angry images are quite similar in appearance, while (for example) Happy and Scream are very dissimilar.

Thus, both for varying lighting (Fig. 15a) and expression (Fig. 15b), training sets with wider variation in the images allow greater robustness to occlusion. Designing a training set that allows recognition under widely varying conditions does not hinder our algorithm; in fact, it helps it. However, the training set should not contain too many similar images, as in the Neutral+Angry example in Fig. 15b. In the language of signal representation, the training images should form an *incoherent dictionary* [9].

of features used (in our face recognition example, approximately 100 are sufficient to make the difference negligible). Moreover, occlusion and corruption can be handled uniformly and robustly within the same classification framework. One can achieve a striking recognition performance for severely occluded or corrupted images by a simple algorithm with no special engineering.

An intriguing question for future work is whether this framework can be useful for object detection, in addition to recognition. The usefulness of sparsity in detection has been noticed in the work in [61] and more recently explored in [62]. We believe that the full potential of sparsity in robust object detection and recognition together is yet to be uncovered. From a practical standpoint, it would also be useful to extend the algorithm to less constrained conditions, especially variations in object pose. Robustness to occlusion allows the algorithm to tolerate small pose variation or misalignment. Furthermore, in the supplementary appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2008.79>, we discuss our algorithm's ability to adapt to nonlinear training distributions. However, the number of training samples required to directly represent the distribution of face images under varying pose may be prohibitively large. Extrapolation in pose, e.g., using only frontal training images, will require integrating feature matching techniques or nonlinear deformation models into the computation of the sparse representation of the test image. Doing so, in a principled manner, it remains an important direction for future work.

## 5 CONCLUSIONS AND DISCUSSIONS

In this paper, we have contended both theoretically and experimentally that exploiting sparsity is critical for the high-performance classification of high-dimensional data such as face images. With sparsity properly harnessed, the choice of features becomes less important than the number

## ACKNOWLEDGMENTS

The authors would like to thank Dr. Harry Shum, Dr. Xiaoou Tang and many others at the Microsoft Research, Asia, for helpful and informative discussions on face recognition, during their visit there in Fall 2006. They also

thank Professor Harm Derksen and Prof. Michael Wakin of the University of Michigan, Professor Robert Fossum and Yoav Sharon of the University of Illinois for the advice and discussions on polytope geometry and sparse representation. This work was partially supported by the Grants ARO MURI W911NF-06-1-0076, US National Science Foundation (NSF) CAREER IIS-0347456, NSF CRS-EHS-0509151, NSF CCF-TF-0514955, ONR YIP N00014-05-1-0633, NSF ECCS07-01676, and NSF IIS 07-03756.

## REFERENCES

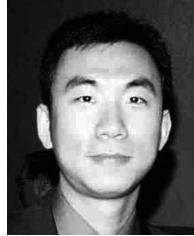
- [1] J. Rissanen, "Modeling by Shortest Data Description," *Automatica*, vol. 14, pp. 465-471, 1978.
- [2] M. Hansen and B. Yu, "Model Selection and the Minimum Description Length Principle," *J. Am. Statistical Assoc.*, vol. 96, pp. 746-774, 2001.
- [3] A. d'Aspremont, L.E. Ghaoui, M. Jordan, and G. Lanckriet, "A Direct Formulation of Sparse PCA Using Semidefinite Programming," *SIAM Rev.*, vol. 49, pp. 434-448, 2007.
- [4] K. Huang and S. Aviyente, "Sparse Representation for Signal Classification," *Neural Information Processing Systems*, 2006.
- [5] V. Vapnik, *The Nature of Statistical Learning Theory*. Springer, 2000.
- [6] T. Cover, "Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition," *IEEE Trans. Electronic Computers*, vol. 14, no. 3, pp. 326-334, 1965.
- [7] B. Olshausen and D. Field, "Sparse Coding with an Overcomplete Basis Set: A Strategy Employed by V1?" *Vision Research*, vol. 37, pp. 3311-3325, 1997.
- [8] T. Serre, "Learning a Dictionary of Shape-Components in Visual Cortex: Comparison with Neurons, Humans and Machines," PhD dissertation, MIT, 2006.
- [9] D. Donoho, "For Most Large Underdetermined Systems of Linear Equations the Minimal  $\ell_1$ -Norm Solution Is Also the Sparsest Solution," *Comm. Pure and Applied Math.*, vol. 59, no. 6, pp. 797-829, 2006.
- [10] E. Candès, J. Romberg, and T. Tao, "Stable Signal Recovery from Incomplete and Inaccurate Measurements," *Comm. Pure and Applied Math.*, vol. 59, no. 8, pp. 1207-1223, 2006.
- [11] E. Candès and T. Tao, "Near-Optimal Signal Recovery from Random Projections: Universal Encoding Strategies?" *IEEE Trans. Information Theory*, vol. 52, no. 12, pp. 5406-5425, 2006.
- [12] P. Zhao and B. Yu, "On Model Selection Consistency of Lasso," *J. Machine Learning Research*, no. 7, pp. 2541-2567, 2006.
- [13] E. Amaldi and V. Kann, "On the Approximability of Minimizing Nonzero Variables or Unsatisfied Relations in Linear Systems," *Theoretical Computer Science*, vol. 209, pp. 237-260, 1998.
- [14] R. Tibshirani, "Regression Shrinkage and Selection via the LASSO," *J. Royal Statistical Soc. B*, vol. 58, no. 1, pp. 267-288, 1996.
- [15] E. Candès, "Compressive Sampling," *Proc. Int'l Congress of Mathematicians*, 2006.
- [16] A. Martinez, "Recognizing Imprecisely Localized, Partially Occluded, and Expression Variant Faces from a Single Sample per Class," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 24, no. 6, pp. 748-763, June 2002.
- [17] B. Park, K. Lee, and S. Lee, "Face Recognition Using Face-ARG Matching," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1982-1988, Dec. 2005.
- [18] R. Duda, P. Hart, and D. Stork, *Pattern Classification*, second ed. John Wiley & Sons, 2001.
- [19] J. Ho, M. Yang, J. Lim, K. Lee, and D. Kriegman, "Clustering Appearances of Objects under Varying Illumination Conditions," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pp. 11-18, 2003.
- [20] S. Li and J. Lu, "Face Recognition Using the Nearest Feature Line Method," *IEEE Trans. Neural Networks*, vol. 10, no. 2, pp. 439-443, 1999.
- [21] P. Sinha, B. Balas, Y. Ostrovsky, and R. Russell, "Face Recognition by Humans: Nineteen Results All Computer Vision Researchers Should Know about," *Proc. IEEE*, vol. 94, no. 11, pp. 1948-1962, 2006.
- [22] W. Zhao, R. Chellappa, J. Phillips, and A. Rosenfeld, "Face Recognition: A Literature Survey," *ACM Computing Surveys*, pp. 399-458, 2003.
- [23] M. Turk and A. Pentland, "Eigenfaces for Recognition," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 1991.
- [24] P. Belhumeur, J. Hespanda, and D. Kriegman, "Eigenfaces versus Fisherfaces: Recognition Using Class Specific Linear Projection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 711-720, July 1997.
- [25] X. He, S. Yan, Y. Hu, P. Niyogi, and H. Zhang, "Face Recognition Using Laplacianfaces," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 3, pp. 328-340, Mar. 2005.
- [26] J. Kim, J. Choi, J. Yi, and M. Turk, "Effective Representation Using ICA for Face Recognition Robust to Local Distortion and Partial Occlusion," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1977-1981, Dec. 2005.
- [27] S. Li, X. Hou, H. Zhang, and Q. Cheng, "Learning Spatially Localized, Parts-Based Representation," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pp. 1-6, 2001.
- [28] A. Leonardis and H. Bischof, "Robust Recognition Using Eigenimages," *Computer Vision and Image Understanding*, vol. 78, no. 1, pp. 99-118, 2000.
- [29] F. Sanja, D. Skocaj, and A. Leonardis, "Combining Reconstructive and Discriminative Subspace Methods for Robust Classification and Regression by Subsampling," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 3, Mar. 2006.
- [30] R. Basri and D. Jacobs, "Lambertian Reflection and Linear Subspaces," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 3, pp. 218-233, 2003.
- [31] H. Wang, S. Li, and Y. Wang, "Generalized Quotient Image," *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pp. 498-505, 2004.
- [32] K. Lee, J. Ho, and D. Kriegman, "Acquiring Linear Subspaces for Face Recognition under Variable Lighting," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 5, pp. 684-698, May 2005.
- [33] D. Donoho and M. Elad, "Optimal Sparse Representation in General (Nonorthogonal) Dictionaries via  $\ell^1$  Minimization," *Proc. Nat'l Academy of Sciences*, pp. 2197-2202, Mar. 2003.
- [34] S. Chen, D. Donoho, and M. Saunders, "Atomic Decomposition by Basis Pursuit," *SIAM Rev.*, vol. 43, no. 1, pp. 129-159, 2001.
- [35] D. Donoho and Y. Tsaig, "Fast Solution of  $\ell^1$ -Norm Minimization Problems when the Solution May Be Sparse," preprint, <http://www.stanford.edu/~tsaig/research.html>, 2006.
- [36] D. Donoho, "Neighboring Polytopes and Sparse Solution of Underdetermined Linear Equations," Technical Report 2005-4, Dept. of Statistics, Stanford Univ., 2005.
- [37] Y. Sharon, J. Wright, and Y. Ma, "Computation and Relaxation of Conditions for Equivalence between  $\ell^1$  and  $\ell^0$  Minimization," CSL Technical Report UILU-ENG-07-2208, Univ. of Illinois, Urbana-Champaign, 2007.
- [38] D. Donoho, "For Most Large Underdetermined Systems of Linear Equations the Minimal  $\ell^1$ -Norm Near Solution Approximates the Sparsest Solution," *Comm. Pure and Applied Math.*, vol. 59, no. 10, 907-934, 2006.
- [39] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge Univ. Press, 2004.
- [40] E. Candès and J. Romberg, " $\ell^1$ -Magic: Recovery of Sparse Signals via Convex Programming," <http://www.acm.caltech.edu/l1magic/>, 2005.
- [41] M. Savvides, R. Abiantun, J. Heo, S. Park, C. Xie, and B. Vijayakumar, "Partial and Holistic Face Recognition on FRGC-II Data Using Support Vector Machine Kernel Correlation Feature Analysis," *Proc. Conf. Computer Vision and Pattern Recognition Workshop (CVPR)*, 2006.
- [42] C. Liu, "Capitalize on Dimensionality Increasing Techniques for Improving Face Recognition Grand Challenge Performance," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 5, pp. 725-737, May 2006.
- [43] P. Phillips, W. Scruggs, A. O'Tools, P. Flynn, K. Bowyer, C. Schott, and M. Sharpe, "FRVT 2006 and ICE 2006 Large-Scale Results," Technical Report NISTIR 7408, NIST, 2007.
- [44] D. Donoho and J. Tanner, "Counting Faces of Randomly Projected Polytopes When the Projection Radically Lowers Dimension," preprint, <http://www.math.utah.edu/~tanner/>, 2007.
- [45] H. Rauhut, K. Schnass, and P. Vanderghenst, "Compressed Sensing and Redundant Dictionaries," to appear in *IEEE Trans. Information Theory*, 2007.
- [46] D. Donoho, "High-Dimensional Data Analysis: The Curses and Blessings of Dimensionality," *AMS Math Challenges Lecture*, 2000.

- [47] S. Kaski, "Dimensionality Reduction by Random Mapping," *Proc. IEEE Int'l Joint Conf. Neural Networks*, vol. 1, pp. 413-418, 1998.
- [48] D. Achlioptas, "Database-Friendly Random Projections," *Proc. ACM Symp. Principles of Database Systems*, pp. 274-281, 2001.
- [49] E. Bingham and H. Mannila, "Random Projection in Dimensionality Reduction: Applications to Image and Text Data," *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 245-250, 2001.
- [50] R. Baraniuk and M. Wakin, "Random Projections of Smooth Manifolds," *Foundations of Computational Math.*, 2007.
- [51] R. Baraniuk, M. Davenport, R. de Vore, and M. Wakin, "The Johnson-Lindenstrauss Lemma Meets Compressed Sensing," *Constructive Approximation*, 2007.
- [52] F. Macwilliams and N. Sloane, *The Theory of Error-Correcting Codes*. North Holland Publishing Co., 1981.
- [53] J. Kim, J. Choi, J. Yi, and M. Turk, "Effective Representation Using ICA for Face Recognition Robust to Local Distortion and Partial Occlusion," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1977-1981, Dec. 2005.
- [54] S. Li, X. Hou, H. Zhang, and Q. Cheng, "Learning Spatially Localized, Parts-Based Representation," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 1-6, 2001.
- [55] T. Ahonen, A. Hadid, and M. Pietikainen, "Face Description with Local Binary Patterns: Application to Face Recognition," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 12, pp. 2037-2041, Dec. 2006.
- [56] M. Lades, J. Vorbruggen, J. Buhmann, J. Lange, C. von der Malsburg, R. Wurtz, and W. Konen, "Distortion Invariant Object Recognition in the Dynamic Link Architecture," *IEEE Trans. Computers*, vol. 42, pp. 300-311, 1993.
- [57] A. Pentland, B. Moghaddam, and T. Starner, "View-Based and Modular Eigenspaces for Face Recognition," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 1994.
- [58] A. Georghiades, P. Belhumeur, and D. Kriegman, "From Few to Many: Illumination Cone Models for Face Recognition under Variable Lighting and Pose," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 6, pp. 643-660, June 2001.
- [59] K. Lee, J. Ho, and D. Kriegman, "Acquiring Linear Subspaces for Face Recognition under Variable Lighting," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 5, pp. 684-698, 2005.
- [60] A. Martinez and R. Benavente, "The AR Face Database," CVC Technical Report 24, 1998.
- [61] D. Geiger, T. Liu, and M. Donahue, "Sparse Representations for Image Decompositions," *Int'l J. Computer Vision*, vol. 33, no. 2, 1999.
- [62] R. Zass and A. Shashua, "Nonnegative Sparse PCA," *Proc. Neural Information and Processing Systems*, 2006.



**John Wright** received the BS degree in computer engineering and the MS degree in electrical engineering from the University of Illinois, Urbana-Champaign. He is currently a PhD candidate in the Decision and Control Group, University of Illinois. His research interests included automatic face and object recognition, sparse signal representation, and minimum description length techniques in supervised and unsupervised learning and has published more than 10 papers on these subjects. He has been a recipient of several awards and fellowships, including the UIUC ECE Distinguished Fellowship and a Carver Fellowship. Most recently, in 2008, he received a Microsoft Research Fellowship, sponsored by Microsoft Live Labs. He is a student member of the IEEE.

than 10 papers on these subjects. He has been a recipient of several awards and fellowships, including the UIUC ECE Distinguished Fellowship and a Carver Fellowship. Most recently, in 2008, he received a Microsoft Research Fellowship, sponsored by Microsoft Live Labs. He is a student member of the IEEE.



**Allen Y. Yang** received the bachelor's degree in computer science from the University of Science and Technology of China in 2001 and the PhD degree in electrical and computer engineering from the University of Illinois, Urbana-Champaign, in 2006. He is a postdoctoral researcher in the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. His primary research is on pattern analysis of geometric or statistical models in very high-dimensional data space and applications in motion segmentation, image segmentation, face recognition, and signal processing in heterogeneous sensor networks. He is the coauthor of five journal papers and more than 10 conference proceedings. He is also the coinventor of two US patent applications. He is a member of the IEEE.



**Arvind Ganesh** received the bachelor's and master's degrees in electrical engineering from the Indian Institute of Technology, Madras, India, in 2006. He is currently working toward the PhD degree in electrical engineering at the University of Illinois, Urbana-Champaign. His research interests include computer vision, machine learning, and signal processing. He is a student member of the IEEE.



**S. Shankar Sastry** received the PhD degree from the University of California, Berkeley, in 1981. He was on the faculty of MIT as an assistant professor from 1980 to 1982 and Harvard University as a chaired Gordon McKay professor in 1994. He served as the chairman of the Department of Electrical Engineering and Computer Sciences, University of California (UC), Berkeley, from 2001 to 2004. He served as the director of the Information Technology Office, DARPA, in 2000. He is currently the Roy W. Carlson professor of electrical engineering and computer science, bioengineering and mechanical engineering, as well as the dean of the College of Engineering, UC Berkeley. He also serves as the director of the Blum Center for Developing Economies. He is the coauthor of more than 300 technical papers and nine books. He received numerous awards, including the President of India Gold Medal in 1977, an M.A. (honoris causa) from Harvard in 1994, fellow of the IEEE in 1994, the distinguished Alumnus Award of the Indian Institute of Technology in 1999, the David Marr prize for the best paper at the Int'l Conference in Computer Vision in 1999, and the Ragazzini Award for Excellence in Education by the American Control Council in 2005. He is a member of the National Academy of Engineering and the American Academy of Arts and Sciences. He is on the Air Force Science Board and is the chairman of the Board of the International Computer Science Institute. He is also a member of the boards of the Federation of American Scientists and Embedded Systems Consortium for Hybrid and Embedded Research (ESCHER).



**Yi Ma** received the bachelors' degrees in automation and applied mathematics from Tsinghua University, Beijing, China, in 1995, the MS degree in electrical engineering and computer science (EECS) in 1997, the MA degree in mathematics in 2000, and the PhD degree in EECS in 2000 from the University of California, Berkeley. Since 2000, he has been on the faculty of the Electrical and Computer Engineering Department, University of Illinois, Urbana-Champaign, where he is currently the rank of associate professor. His main research interests include systems theory and computer vision. He was the recipient of the David Marr Best Paper Prize at the International Conference on Computer Vision in 1999 and Honorable Mention for the Longuet-Higgins Best Paper Award at the European Conference on Computer Vision in 2004. He received the CAREER Award from the National Science Foundation in 2004 and the Young Investigator Program Award from the Office of Naval Research in 2005. He is a senior member of the IEEE and a member of the ACM.