



Weighted Huber constrained sparse face recognition

Dajiang Lei^{1,2} · Zhijie Jiang¹ · Yu Wu²

Received: 25 October 2018 / Accepted: 8 January 2019 / Published online: 21 January 2019
© Springer-Verlag London Ltd., part of Springer Nature 2019

Abstract

Recently sparse coding based on regression analysis has been widely used in face recognition research. Most existing regression methods add an extra constraint factor to the coding residual to make the fidelity term in the l_2 loss approach the Gaussian or Laplace distribution. But the essence of these methods is that only the fidelity term of l_1 loss or l_2 loss is used. In this paper, weighted Huber constrained sparse coding (WHCSC) is used to study the robustness of face recognition in occluded environments, and alternating direction method of multipliers is used to solve the problem of model minimization. In WHCSC, we propose a sparse coding with weight learning and use Huber loss to determine whether the fidelity is a l_2 loss or l_1 loss. For the WHCSC model, the two kinds of classification modes and the two kinds of weight coefficients are further studied for the intra-class difference and the inter-class difference in the face image classification. Through a large number of experiments on a public face database, WHCSC shows strong robustness in face occlusion, corrosion and illumination changes comparing to the state-of-the-art methods.

Keywords Sparse coding · Face recognition · Robustness · ADMM

1 Introduction

In recent years, face recognition is still a hot research topic [1]. On the one hand, it is great potential for use; on the other hand, it reveals how machine learning can make feature selection and classification on complete images [2, 3]. The advantage of face recognition lies in its naturalness and the characteristics that are not perceived by the tested individual [4]. First, naturalness means that the recognition method is the same as the biological characteristics used by human (or even other organisms) for individual recognition. For example, in face recognition, humans beings also distinguish and confirm identity by

observing face. Second, unobtrusive characteristics are also important for a method of identification, which makes it less objectionable and because it is not easy to attract people's attention, it is not easy to be deceived. Face recognition uses visible light to acquire face image information. This is different from fingerprint recognition or iris recognition, which requires the use of an electronic pressure sensor to capture fingerprints, or the use of infrared to acquire iris images. Fingerprint recognition or iris recognition is easily perceived and thus more likely deceived by camouflage.

Face recognition is considered as one of the most difficult research topics in the field of biometrics and even artificial intelligence. On the one hand, this difficulty comes from the characteristics of human biological characteristics. First, the similarity of the face: the structure and appearance of the faces between different individuals are very similar. This similarity is not conducive to the use of human face to distinguish between human beings. Second, the variability of the face: face shape is very unstable, human complex facial expression changes, but also in different angles of view, face the visual image is also very different. Finally, the difference in face: different genetic makeup makes each person's face always different. On the other hand, the external noise changes. For example,

✉ Dajiang Lei
leidj@cqupt.edu.cn

Zhijie Jiang
S160231034@stu.cqupt.edu.cn

Yu Wu
wuyu@cqupt.edu.cn

¹ College of Computer, Chongqing University of Posts and Telecommunications, Chongqing 400065, China

² Institute of Web Intelligence, Chongqing University of Posts and Telecommunications, Chongqing 400065, China

expression changes, lighting conditions, true camouflage, continuous occlusion, pixel corrosion, etc., will reduce the valuable information of the face image and interfere with the recognition of the face [5].

Recently, the method based on regression analysis has attracted the widespread interest of researchers. The linear regression classifier (LRC) proposed by Naseem et al. [6] represents the query image by linear combination of dictionary atoms. Wright et al. [7] proposed a sparse coding classification algorithm (SRC) to identify the real camouflage and pixel erosion of human face images. The SRC uses a sparse linear combination of dictionary atoms to represent the query image. LRC and SRC cannot achieve the desired performance when the dictionary is not enough. Zhang et al. [8] considered that cooperative mechanisms are more important than sparse constraints. They proposed a collaborative representation classifier (CRC) based on l_2 norm constraints and further proposed its robust version (RCRC). Yang et al. [9] further proposed a matrix regression (NMR) classification framework based on kernel regularization of l_2 norm to obtain a better recognition rate in occlusion and lighting changes. Zhong et al. [10] considered that the balance between SRC and CRC can be adjusted through iteration and a classifier (LHC) of $l_{1/2}$ regularization and ITR iterative mechanism is proposed. Zheng et al. [11] attempted to obtain a more general classifier (IRGSC) through adaptive feature weight learning and adaptive distance-weighted learning. Lin et al. [12] propose a robust, discriminative and comprehensive dictionary learning (RDCDL) method, in which a robust dictionary is learned from comprehensive training sample diversities generated by extracting and generating facial variations.

Although these classifiers have made great progress, due to the complicated changes of occlusion, two types of changes of face image are still not well overcome. The first type of change is called inter-class difference. The inter-class difference should be amplified as a standard to distinguish between individuals. The second type of change is called intra-class difference. It should be eliminated because they can represent the same individual. For face images, intra-class difference interference is often greater than the inter-class difference, so it becomes very difficult to distinguish the individuals by the inter-class difference under the intra-class difference. Two types of changes are one of the biggest obstacles that face recognition technology is widely used and need to be solved urgently.

We propose a new scheme called weighted Huber constrained sparse coding (WHCSC) and establish a robust weighted regression model with sparse constraints. WHCSC seeks the problem with the maximum a posteriori (MAP) of sparse coding and is robust against noise values (such as occlusion, corrosion, illumination changes). The

experiment uses a representative face database. The experimental results show that WHCSC has obvious advantages in dealing with facial occlusion, corrosion, and camouflage.

The main contributions of this article are summarized as follows:

1. Propose a more efficient taxonomy. On the one hand, the intra-class difference is reduced; on the other hand, the inter-class interference is effectively avoided when the coding coefficients of the query sample and the training sample are calculated.
2. The weighted method is adopted to reduce the influence of noise. At the same time, two kinds of exponential form of weight are researched to further expand the effect of weight vector, increase the inter-class difference and improve the recognition rate.
3. Utilize the robustness of Huber function to reduce outlier interference and solve the l_1 norm minimization problem by using alternating direction method of multipliers (ADMM) [13].

The rest of the paper is organized as follows: Sect. 2 introduces sparse robust coding. Section 3 introduces the WHCSC and its contributions. Section 4 analyzes the computational complexity of WHCSC. Section 5 analyzes the convergence of WHCSC. Section 6 tests WHCSC performance using a published face dataset. Finally, Sect. 7 summarizes WHCSC.

2 Sparse robust coding based on regression classifier

2.1 General classification framework based on regression analysis

In the general classification problem, the training samples are expressed as a dictionary matrix $X = [X_1, X_2, \dots, X_c] \in \mathbf{R}^{m \times n}$, and c is the sample category. $X_i = [X_{i1}, X_{i2}, \dots, X_{in_i}] \in \mathbf{R}^{m \times n_i}$ ($i = 1, 2, \dots, c$) is the sample subset of each category of the sample corpus X . n_i is the number of training samples of class i , $n = \sum_{i=1}^c n_i$ is the total number of samples. In the regression, the training sample X linearly represents the query sample

$$y = X_1\theta_1 + X_2\theta_2 + \dots + X_c\theta_c \\ = X_{11}\theta_{11} + X_{12}\theta_{12} + \dots + X_{cn_c}\theta_{cn_c} = X\theta, \quad (1)$$

where $\theta = [\theta_{11}, \theta_{12}, \dots, \theta_{cn_c}]^T \in \mathbf{R}^n$ is the coding coefficient of the query sample to be determined on the training sample.

The regression-based classification is to determine the class of the query sample $y \in \mathbf{R}^m$ in a given training sample. By computing the residuals $e_i = y - X_i\theta_i$ in the

query sample and each category, the category of the smallest e_i is regarded as the category of the query sample.

2.2 Sparse coding

Sparse coding is an artificial neural network method that simulates the simple cell receptivity field in the primary visual cortex V1 of the mammalian visual system and has been widely used in image processing and natural language [14, 15]. Some human visual studies suggest that many neurons in the visual pathway are selective for a variety of specific stimuli in lower- and intermediate-level human vision, such as color, texture, orientation, size. [16, 17]. Given the sparseness of the input image given by these neurons, it can be efficiently computed by convex optimization. Due to the difficulty in solving the l_0 norm minimization, the l_1 norm is usually used as the nearest solution to the l_0 norm minimization problem. In general, the problem of sparse coding can be expressed as

$$\min_{\theta} \|y - X\theta\|_2^2 + \lambda\|\alpha\|_1 \quad \text{s.t. } \alpha = \theta, \tag{2}$$

where λ is the penalty coefficient for the l_1 norm. The essence of formula (2) is the least squares estimation of sparse constraints when the residuals follow a Gaussian distribution. When the residuals follow the Laplace distribution, the sparse coding problem is

$$\min_{\theta} \|y - X\theta\|_1 + \lambda\|\alpha\|_1 \quad \text{s.t. } \alpha = \theta. \tag{3}$$

Sparse coding can capture high-order correlation structures in an image and represent the signal with as few atoms as possible in a given overcomplete dictionary [18]. However, there are mainly two problems with this model. The first one is whether the regularized l_1 norm constraint $\|\alpha\|_1$ is good enough to make the signal sufficiently sparse. The second one is whether the fidelity term ($\|y - X\theta\|_2^2$ or $\|y - X\theta\|_1$) is sufficiently effective to describe the fidelity of the signal, especially when the signal has noise or abnormal values.

Improve the first problem by modifying sparse constraints. For example, Liu et al. [19] added a nonnegative constraint on sparse coefficient α . Gao et al. [20] introduced Laplace coefficients in sparse coding. Wang et al. [21] used weighted l_2 norm for sparse constraints. In addition, Ramirez et al. [22] proposed a generic sparse modeling framework to design sparse regularization terms.

For the second problem, defining the fidelity terms using the l_2 or l_1 norm from the perspective of the maximum a posteriori probability (MAP) actually assumes that the encoded residuals follow a Gaussian or Laplace distribution. However, in practice, it may not be very good to follow a certain distribution of a single hypothetical residual, especially when occlusion, camouflage or

corruption occurs in facial images. Therefore, a fidelity item that uses a single l_2 or l_1 norm in a sparse coding model may not be robust in these cases.

2.3 Sparse robust coding

It can be observed from (c) in Fig. 1 that when the encoding residual approaches 0, the encoding residual of the l_2 norm is smaller, and when it is far from 0, the l_1 norm is smaller.

In practice, in a large number of training samples, it will naturally contain more or less some outliers. In linear coding, it is assumed that the sum of the residuals of the training samples and the query samples is $\sum_{i=1}^m e_i$, and the outliers have a great contribution to $\sum_{i=1}^m e_i$. Therefore, to some extent reduce the encoding residuals of outliers, will be greatly reduced $\sum_{i=1}^m e_i$. For example, in Fig. 1 (c), l_2 loss and l_1 loss show two different coding residuals. Therefore, in order to reduce the impact of outliers, it is important to query for different pixels using different fidelities ($\|y - X\theta\|_2^2$ or $\|y - X\theta\|_1$).

In the statistical learning perspective, the Huber loss function is a loss function of robust regression, which is insensitive to outliers compared to mean square error and is often used for classification problems. Huber loss function is expressed as

$$g(z) = \begin{cases} |z|^2/2 & |z| \leq \eta \\ \eta|z| - \eta^2/2 & |z| > \eta \end{cases}, \tag{4}$$

where z is the residual and η is the Huber threshold.

l_2 loss and l_1 loss are mixed in the Huber loss (Fig. 2). If the absolute value of the residual $|z|$ is smaller than the threshold value η (that is, the normal value), the fidelity of formula (4) uses l_2 loss. If the absolute value of the residual the value $|z|$ is greater than the threshold η , and the fidelity of formula (4) uses l_1 loss. For smooth connection with l_2 loss, the constant $\eta^2/2$ is subtracted from l_1 loss. The Huber loss balances the validity and robustness through the optimal combination of l_2 loss and l_1 loss.

In order to improve the robustness and validity of sparse coding, a sparse Huber (SH) model is designed according to Huber loss mentioned above.

Byod explains in Sect. 6 of the article [13] that the Huber function corresponds to the standard form $\min f(\theta) + g(z)$ of the ADMM model and can be expressed as:

$$\min_{\theta} g(z) \quad \text{s.t. } z = X\theta - y, \tag{5}$$

$$\text{where } f(\theta) = 0, g(z) = \begin{cases} \frac{1}{2}\|z\|_2^2 & |z| \leq \eta \\ \eta\|z\|_1 - \frac{1}{2}\eta^2|z| & |z| > \eta \end{cases} \quad \text{and is constrained by } z = y - X\theta.$$

Fig. 1 **a** and **b** are two pictures of the same person. **c** is an l_2 loss and l_1 loss coded residual image of the coded residual of **(a)** and **(b)** when the residual threshold is 10

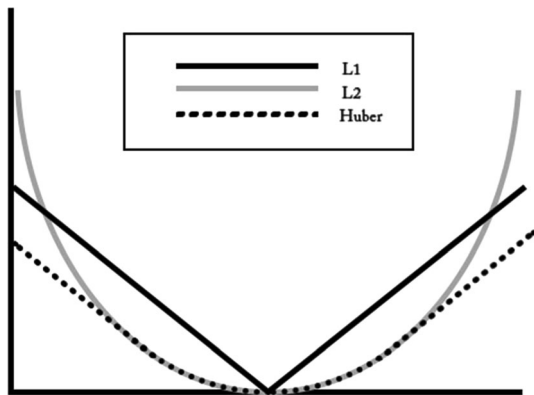
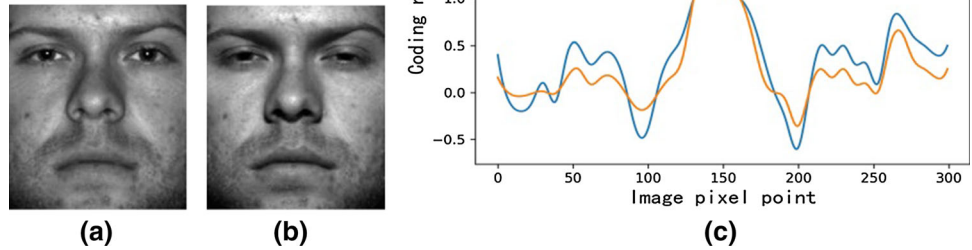


Fig. 2 Huber loss function diagram

SH can be expressed as

$$\min_{\theta} g(z) + \lambda \alpha_1 \quad \text{s.t. } z = X\theta - y, \alpha = \theta, \quad (6)$$

where $\lambda \|\alpha\|_1$ is an l_1 norm penalty term with $\alpha = \theta$ constraints. In a certain range, the larger the value of λ , the more sparse θ .

3 Weighted Huber constrained sparse coding

3.1 Weighted Huber constrained sparse coding model

To further reduce the effects of noise or outliers in the training samples, we design a weight for the training samples so that outliers are given a low weight value. In RSRC, an effective weight vector is proposed to convert the minimization problem into an iteratively reweighted sparse coding problem.

With reference to the weight vector in RSRC [23], combined with the above-mentioned sparse Huber model (SH), this paper proposes a weighted Huber constrained sparse coding model (WHCSC). The WHCSC model is

essentially a maximum likelihood estimation (MLE) problem. The weight vector and the SH model are jointly used to reduce the noise interference, and ADMM is used to solve the l_1 norm minimization problem. A large number of experiments conducted in open face database show that WHCSC has good classification effect, especially when the facial image has complex changes such as occlusion, corrosion, light changes.

WHCSC can be expressed as

$$\min_{\theta} g(z) + \lambda \|\alpha\|_1 \quad \text{s.t. } z = w \odot (X\theta - y), \alpha = \theta, \quad (7)$$

$$\text{where } f(\theta) = 0, g(z) = \begin{cases} \frac{1}{2} \|z\|_2^2 & |z_k| \leq \eta w_k \\ \eta \|w \odot z\|_1 - \frac{1}{2} \eta^2 w^T w & |z_k| > \eta w_k \end{cases}$$

$k = 1, 2, \dots, m$, w is the sample weight. η is the residual threshold constant. There are different methods to determine the threshold of η in many papers. In this paper, we propose a combined threshold of weight, that is ηw , where η is a constant. ηw makes the threshold value more in line with the distribution of training samples with weight w constraint. $a \odot b$ represents the multiplication of the corresponding elements of a and b .

Weight $w = [w_1, w_2, \dots, w_m] \in \mathbf{R}^{m \times 1}$. w_m is the weight of number m in training sample $X \in \mathbf{R}^{m \times n}$, e_m is the residual with number m , w_m is set to the following sigmoid function

$$w_m(e_m) = \frac{1}{1 + \exp\left(-q \left(\frac{\delta - e_m^2}{\delta}\right)\right)}, \quad (8)$$

where δ is the residual threshold. $\delta - e_m^2$ represents the distance between the residual and the residual threshold, and $\frac{\delta - e_m^2}{\delta}$ unifies the dimension of this distance. q affects the penalty rate of weights and makes the distribution of weights smoother. The sigmoid function can constrain the weights value between $[0, 1]$. Therefore, when the residual is greater than δ , the weights is less than 0.5; when equal to

δ , the weights is equal to 0.5; when less than δ , the weights is greater than 0.5.

Let $\Psi = [e_1^2, e_2^2, \dots, e_m^2]$, and then ranking Ψ to get Ψ_a . Let $k = \lfloor \tau m \rfloor$, where $\tau \in (0, 1]$, $\lfloor \tau m \rfloor$ is an integer less than $\lfloor \tau m \rfloor$, then δ can be expressed as

$$\delta = \Psi_a(k). \tag{9}$$

For ease of calculation, formula (8) is organized to get

$$w_m(e_m) = \frac{\exp(-\mu e_m^2 + \mu\delta)}{1 + \exp(-\mu e_m^2 + \mu\delta)}, \tag{10}$$

where the parameter $\mu = \frac{q}{\delta}$.

Compared with the model in formulas (6), (7) has the following advantages. Outliers (usual pixels with large residuals) are adaptively assigned a low weight to reduce their impact on regression estimates, which can greatly reduce the sensitivity to outliers. And formula (8) limits the weights between [0,1] using a sigmoid function and avoids the almost infinite weight value of pixels with very small residuals, which improves the stability of the encoding process. The important parameters q and τ will be analyzed in conjunction with the experiment in Sect. 6.6.

3.2 WHCSC’s contribution

The purpose of the linear expression-based classification is to obtain the smallest encoding residual by linear expression with the optimal encoding coefficient θ , thereby distinguishing the category to which the test image belongs. Definition

$$y_i = F_i(\mathbf{X}) = [\mathbf{X}_1\theta_1 + \mathbf{X}_2\theta_2 + \dots + \mathbf{X}_i\theta_i + \dots + \mathbf{X}_c\theta_c],$$

where $y_i = F_i(\mathbf{X})$ represents the linear expression of the sample set for the i th test sample. Due to the variability of the human face, two face images generated by the same person at different times do not appear to be identical, resulting in an intra-class difference, that is, $y_i - \mathbf{X}_i\theta_i > 0$. Similarly, the inter-class difference is the difference between different people, that is, $y_i - \mathbf{X}_j\theta_j > 0 (j \neq i)$. Sparse coding has the function of feature selection. Its purpose is to select the training samples that are most similar to the test samples to be linearly combined into test samples [23]. First, we prefer to select samples belonging to the same class to linearly combine test samples and exclude interference from other classes of samples. This makes the coding coefficients (θ_j) of samples of different categories small enough. Second, we also prefer to select the same type of training samples that have less interference with the test samples for the same type of samples. On the other hand, in actual tests, the linear expression of test samples in each category will be calculated. Therefore, we hope that the coding residuals of linear expressions in the same category will be small, while the coding residuals of

different categories are large. Section 3.2.1 describes in detail the methods used to reduce intra-class difference and avoid inter-class interference. Section 3.2.2 describes in detail how to increase the inter-class difference.

3.2.1 Reduce the intra-class difference and remove inter-class interference

In Fig. 3, we assume that the query sample belongs to the category i . The residuals of the query samples and the training samples of each category are $e = [e_{i,1}, e_{i,2}, \dots, e_{i,i}, \dots, e_{i,c}] = [(\mathbf{X}_1\theta_1 - y_i), (\mathbf{X}_2\theta_2 - y_i), \dots, (\mathbf{X}_i\theta_i - y_i), \dots, (\mathbf{X}_c\theta_c - y_i)]$, where $e_{i,i}$ represents the coding residuals of the query samples and training samples of the same class, that is, intra-class difference, and $e_{i,j}$ denotes the encoding residuals of the query samples and the training samples of different class, that is, inter-class difference.

In RSRC, the weight w is defined based on the residual of the complete sample $X \in R^{m \times n}$ and the query sample. And all categories of samples use the same weight vector, i.e. $w \odot (\mathbf{X}\theta - y) = [w \odot (\mathbf{X}_1\theta - y), w \odot (\mathbf{X}_2\theta - y), \dots, w \odot (\mathbf{X}_c\theta - y)]$. Here is defined as a classification model I, as shown in Fig. 3a, which shows the use of the complete works samples to define weights, the same type of coding residuals and the distribution of different types of coding residuals. However, in WHCSC, the weight w is based on the residual definition of the sample subset $\mathbf{X}_i \in R^{m \times n_i}$ and the query sample, that is, $w \odot (\mathbf{X}\theta - y) = [w_1 \odot (\mathbf{X}_1\theta - y), w_2 \odot (\mathbf{X}_2\theta - y), \dots, w_c \odot (\mathbf{X}_c\theta - y)]$. Here is defined as classification model II, as shown in Fig. 3b, which shows the distribution of the same class coded residual and the different class coded residuals when the weights defined by the sample subset are used.

Observing Fig. 3a–c, using the residuals of the sample subset and the query sample to define the weights can significantly reduce the coding residuals of the same class, although the different types of coding residuals also decrease, However, the same type and different types of residuals fit straight line is still a clear distinction.

Therefore, in classification model II, the weight of WHCSC can obtain the independent weights that are more suitable for the subset of samples in this category, and then the better coding coefficient θ_i under this weight is obtained, so as to reduce the intra-class difference.

3.2.2 Increase the inter-class difference

First, it is also assumed that $e = [e_{i,1}, e_{i,2}, \dots, e_{i,i}, \dots, e_{i,c}] = [(\mathbf{X}_1\theta_1 - y_i), (\mathbf{X}_2\theta_2 - y_i), \dots, (\mathbf{X}_i\theta_i - y_i), \dots, (\mathbf{X}_c\theta_c - y_i)]$. The difference between the residuals in the different types of training samples and the residuals in

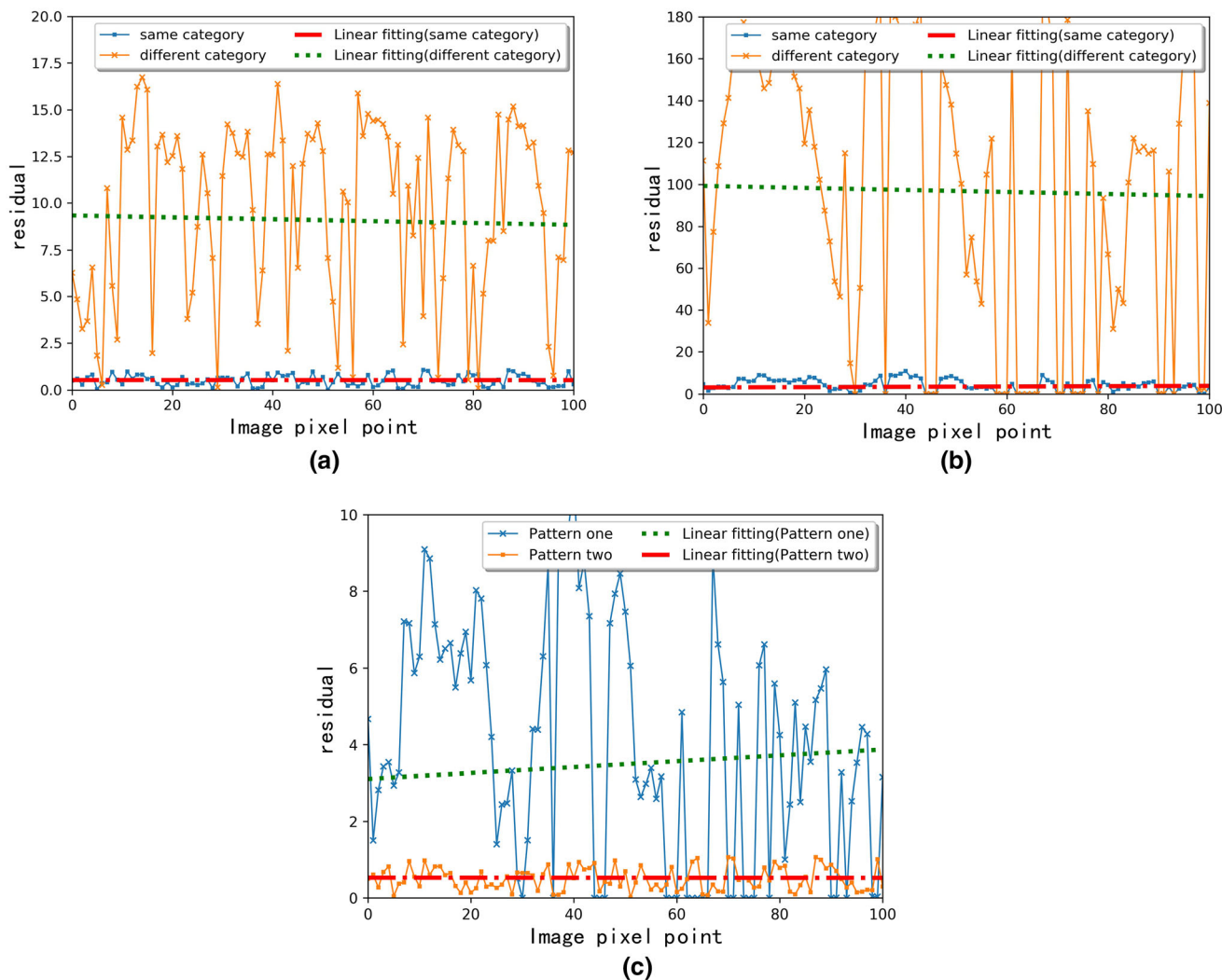


Fig. 3 **a** is a classification mode I. The “different category” curve is the fitted image residual distribution of the query sample and different categories of training samples, and “Linear fitting (different category)” is the corresponding residual distribution fitting straight line. The “same category” curve is the fitted image residual distribution of the query sample and the same category of training samples, and “Linear fitting (same category)” is the corresponding residual distribution fitting straight line. **b** is a classification mode II. The “different category” curve is the fitted image residual distribution of the query sample and the different categories of training samples, and “Linear fitting (different category)” is the corresponding residual distribution fitting straight line. The “same category” curve is the

fitted image residual distribution of the query sample and the same category of training samples, and “Linear fitting (same category)” is the corresponding residual distribution fitting straight line. **c** is a fitted image residual distribution of the query sample and the same category of training samples in the classification mode I and the classification mode II. The “Pattern one” curve and the “Linear fitting (Pattern one)” line respectively correspond to the residual distribution and the residual distribution fitting line in the classification mode I. The “Pattern two” curve and the “Linear fitting (Pattern two)” line respectively correspond to the residual distribution and the residual distribution fitting line in the classification mode II

the category i training samples, i.e. the relative differences between the inter-class difference and intra-class difference, is expressed as

$$\Delta e = \left[\frac{e_{i,1} - e_{i,i}}{e_{i,i}}, \frac{e_{i,2} - e_{i,i}}{e_{i,i}}, \dots, \frac{e_{i,i-1} - e_{i,i}}{e_{i,i}}, \frac{e_{i,i+1} - e_{i,i}}{e_{i,i}}, \dots, \frac{e_{i,c} - e_{i,i}}{e_{i,i}} \right] \in \mathbf{R}^{c-1}. \tag{11}$$

The larger the Δe is, the larger inter-class difference is relative to the intra-class difference, the easier it is to distinguish between the query sample and other types of samples. Conversely, the smaller the Δe is, the smaller inter-class difference is relative to the intra-class difference, the more difficult it is to distinguish between the query sample and other classes of samples.

In RSRC, the weighting effect on the residual is expressed as $w^{\frac{1}{2}} \odot (X\theta - y)$, and the definition of $w^{\frac{1}{2}}$ is 0.5

power exponent weights and the relative difference in residual is $\Delta e_{w^{\frac{1}{2}}}$. However in WHCSC, the weighting effect on the residual is expressed as $w \odot (X\theta - y)$, and the definition of w is 1 power exponent weights, and the relative difference in residual is Δe_w . Figure 4 shows the experimental results of $w^{\frac{1}{2}}$ and w in WHCSC. The ordinate is the distribution of $\Delta e_{w^{\frac{1}{2}}}$ and Δe_w , and the abscissa is the sample type.

It can be observed that the 1 power exponent weight makes the difference between $e_{i,j}$ and $e_{i,i}$ increase, that is to say, the inter-class difference is more different from the intra-class difference, that is, increase the inter-class difference.

3.3 The initial value of the weight

A good initial value will make the algorithm easier to get good performance. In order to initialize the weights, the coding residuals of the query samples should first be estimated. We can set the initial residual as $e = y - y_{mni}$. Because the weight of WHCSC is sub-category calculation, it is reasonable to set y_{mni} as the average of the same pixels of the current training sample subset

$$\begin{aligned}
 y_{mni} &= [m(y_1), m(y_2), \dots, m(y_k)] \\
 &= [m([y_{11}, y_{12}, \dots, y_{1j}]), m([y_{21}, y_{22}, \dots, y_{2j}]), \\
 &\quad \dots, m([y_{k1}, y_{k2}, \dots, y_{kj}]))] \quad k = [1, 2, \dots, m], \\
 j &= [1, 2, \dots, n_i]
 \end{aligned} \tag{12}$$

where $m(x)$ represents the mean of x . For parameters τ and q , usually $\tau = 0.8$ and $q = 1$. In more complex environments, such as occlusion, camouflage, corrosion, you can set smaller τ .

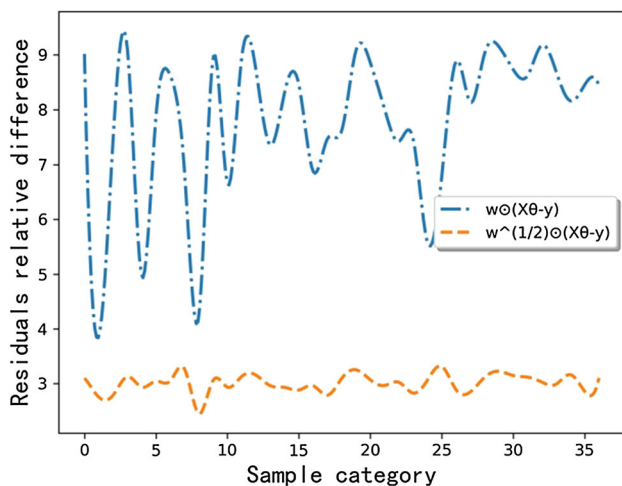


Fig. 4 In two different weighting coefficients, the relative difference between the intra-class difference and the inter-class difference

3.4 WHCSC iteration conditions

In each iteration, the formula (7) will gradually decrease, the lower bound is 0, and WHCSC will gradually converge. WHCSC converges and the iteration terminates when the difference in θ between adjacent iterations is small enough. The termination conditions are as follows

$$\|\theta^t - \theta^{t-1}\|_2^2 < \gamma \tag{13}$$

where γ is a small enough positive number and t is the number of iterations.

3.5 ADMM solves the sub-problem

$$\min_{\theta} g(z) + \lambda \|\alpha\|_1 \quad \text{s.t. } z = w \odot (X\theta - y), \alpha = \theta. \tag{14}$$

The Lagrange expression of a sub-problem is

$$\begin{aligned}
 \mathcal{L}(\theta, z, \alpha, h_z, h_\alpha) &= g(z) + \lambda \|\alpha\|_1 + \langle h_z, w \odot (X\theta - y) - z \rangle \\
 &\quad + \langle h_\alpha, \theta - \alpha \rangle
 \end{aligned} \tag{15}$$

ADMM is an algorithm that aims to fuse the dual variable ascent method’s decomposability and the multiplier method’s upper bound convergence property. In order to increase the robustness of the dual variable ascent method and the strong convex constraint of the relaxation function, introducing the augmented Lagrangian formula

$$\begin{aligned}
 \mathcal{L}_{\rho_1, \rho_2}(\theta, z, \alpha, h_z, h_\alpha) &= g(z) + \lambda \|\alpha\|_1 + \langle h_z, w \odot (X\theta - y) \\
 &\quad - z \rangle + \frac{\rho_1}{2} \|w \odot (X\theta - y) - z\|_2^2 \\
 &\quad + \langle h_\alpha, \theta - \alpha \rangle + \frac{\rho_2}{2} \|\theta - \alpha\|_2^2
 \end{aligned} \tag{16}$$

where ρ_1, ρ_2 is greater than zero. The ADMM iteration is made up of

$$\theta^{k+1} := \arg \min_x \mathcal{L}_{\rho_1, \rho_2}(\theta^k, z^k, \alpha^k, h_z^k, h_\alpha^k) \tag{17}$$

$$z^{k+1} := \arg \min_z \mathcal{L}_{\rho_1}(\theta^{k+1}, z^k, h_z^k) \tag{18}$$

$$\alpha^{k+1} := \arg \min_\alpha \mathcal{L}_{\rho_2}(\theta^{k+1}, \alpha^k, h_\alpha^k) \tag{19}$$

$$h_z^{k+1} := h_z^k + \rho_1 (w \odot (X\theta^{k+1} - y) - z^{k+1}) \tag{20}$$

$$h_\alpha^{k+1} := h_\alpha^k + \rho_2 (\theta^{k+1} - \alpha^{k+1}) \tag{21}$$

Formula (16) is brought into formula (17), (18), (19), (20) and (21), and the iterative step of ADMM is

$$\theta^{k+1} := \arg \min_x \left(\frac{\rho_1}{2} \|w \odot (X\theta^k - y) - z^k + u^k\|_2^2 + \frac{\rho_2}{2} \|\theta^k - \alpha^k + u^k\|_2^2 \right) \tag{22}$$

$$z^{k+1} := \arg \min_z \left(g(z^k) + \frac{\rho_1}{2} \|w \odot (X\theta^{k+1} - y) - z^k + u^k\|_2^2 \right) \quad \alpha^{k+1} = S_{\frac{\rho_2}{2}}(\theta^{k+1} + u^k), \tag{23}$$

$$u_z^{k+1} = u_z^k + w \odot (X\theta^{k+1} - y) - z^{k+1}, \tag{30}$$

$$\alpha^{k+1} := \arg \min_{\alpha} \left(\lambda \|\alpha^k\|_1 + \frac{\rho_2}{2} \|\theta^{k+1} - \alpha^k + u^k\|_2^2 \right) \quad u_{\alpha}^{k+1} = u_{\alpha}^k + \theta^{k+1} - \alpha^{k+1}, \tag{24}$$

$$u_{\alpha}^{k+1} = u_{\alpha}^k + \theta^{k+1} - \alpha^{k+1}, \tag{31}$$

$$u_z^{k+1} := u_z^k + w \odot (X\theta^{k+1} - y) - z^{k+1} \tag{25}$$

where $W = \text{diag}(w) = \text{diag}([w_1, w_2, \dots, w_m])$ and the S operator is defined as

$$u_{\alpha}^{k+1} := u_{\alpha}^k + \theta^{k+1} - \alpha^{k+1} \tag{26}$$

$$S_k(a) = \begin{cases} a - k, & a > k \\ 0, & |a| \leq k \\ a + k, & a < -k \end{cases} \tag{32}$$

where u is an alternative variable for $u = \frac{h}{\rho}$. Solve the formula (22), (23), (24), (25) and (26) to get

$$\theta^{k+1} = (\rho_1 X^T W^T W X + \rho_2)^{-1} [\rho_1 X^T W^T (z^k - u^k + y) + \rho_2 (\alpha^k - u^k)], \tag{27}$$

$$z^{k+1} = \frac{\rho_1}{1 + \rho_1} [w \odot (X\theta^{k+1} - y) + u^k] + S_{\frac{w_i}{\rho_1}} \left[\frac{1}{1 + \rho_1} (w \odot (X\theta^{k+1} - y) + u^k) \right], \tag{28}$$

3.6 Judgment query sample category

The residuals of the query sample in each category are calculated according to the categories θ_i obtained $e = [e_1, e_2, \dots, e_i], i = 1, 2, \dots, c$, where $e_i = y - X_i \theta_i$. The category of the smallest e_i belongs to the category of the query sample.

Algorithm 1 Weighted Huber Constraint Sparse Coding

Input: test sample y , training sample subset X_i , initialization y_{rec}^1 to y_{mn_i} , parameters τ and q , threshold η

Output: θ

①: $i = 1 \in [1, 2, \dots, c]$

②: $t = 1$

③: Calculate the residuals $e_{im}^t = y - y_{rec}^t$

④: The weight is calculated as

$$w_i(e_{im}^t) = \frac{\exp(-\mu^t (e_{im}^t)^2 + \mu^t \delta^t)}{1 + \exp(-\mu^t (e_{im}^t)^2 + \mu^t \delta^t)}$$

⑤: Encoding coefficient update

$$\theta_i^t = \min_{\theta_i} g(z) + \lambda \|\alpha\|_1 \quad s.t. \quad z = w_i^t \odot (X_i \theta_i - y), \alpha = \theta_i$$

the specific solution to sub-problems in section 3.5 introduced.

⑥: Reconstruction $y_{rec}^t = X_i \theta_i^t$, and let $t = t + 1$

⑦: $\theta_i = \theta_i^t$

⑧: Return to step ③ until the conditions are met or the maximum number of iterations is reached.

⑨: After the categories i 's θ_i is updated, let $i = i + 1$ until all the categories θ are updated.

4 Computational complexity analysis

The computational cost of the algorithm is mainly used to update the weight w and the coding coefficient θ . Given that there are m face data sets of one category, and each image size is $n = p \times q$. The face data set has a total of c categories. The number of iterations of algorithm step 2 is denoted as k_1 . The computational complexity of the weight $w \in R^{n \times 1}$ in step 4 is $O(n)$. WX and Wy can be calculated and cached in advance. The computational complexity of θ in formula (27) is $O(nm^2)$, z in formula (28) is $O(nm)$ and u_z in formula (30) is $O(nm)$. Therefore, the computational complexity of the coding coefficient θ in step 5 is $O(k_2nm^2)$, where k_2 is the number of iterations of the ADMM algorithm. In summary, the computational complexity of WHCSC is $O(ck_1(n + k_2nm^2))$ [26, 27]. After many experiments, k_1 and k_2 are usually less than 10.

5 Convergence and convergence rate analysis

Before proofing of convergence, the standard form of the ADMM objective function is given by formula (7) as follows

$$\min f(\theta) + g(z) + l(\alpha) \quad \text{s.t. } z = w \odot (X\theta - y), \alpha = \theta, \tag{33}$$

where $f(\theta) = 0, \quad l(\alpha) = \lambda\alpha_1, \quad g(z) = \begin{cases} \frac{1}{2} \|z\|_2^2 & |z_k| \leq \eta w_k \\ \eta \|w \odot z\|_1 - \frac{1}{2} \eta^2 w^T w & |z_k| > \eta w_k \end{cases}, \quad k = 1, 2, \dots, m.$ The

following are two theorems about the function $f(\theta), g(z), l(\alpha)$.

Theorem 1 *The function of $f(\theta), g(z), l(\alpha)$ is closed, proper, and convex.*

Proof Obviously, $f(\theta) = 0$ must be a closed, proper, and convex function. Since $\lambda > 0$, the norm satisfies the triangle inequality; $l(\alpha)$ is a proper closed convex function. The epigraph of $g(z)$ can be expressed as the following form, i.e.

$$\text{epig} = \{(z, t_z) \in R^m \times R \mid g(z) \leq t_z\}. \tag{34}$$

Obviously the epigraph of $g(z)$ is a non-empty closed convex set. According to the nature of the epigraph, $g(z)$ is a proper closed convex function when epig is a non-empty closed convex function. The iterative step of ADMM algorithm is to solve the optimal solution of each sub-problem. Obviously, the optimal solution of sub-problems

$\theta^{k+1}, z^{k+1}, \alpha^{k+1}$ is feasible. The problem of minimizing $\theta^{k+1}, z^{k+1}, \alpha^{k+1}$ has solution (not necessarily unique). Therefore, $f(\theta), g(z), l(\alpha)$ are proper closed convex functions, and $f(\theta) + g(z) + l(\alpha)$ is also a proper closed convex function. Certificate completed.

Theorem 2 *The unaugmented Lagrangian*

$$\mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha) = g(z) + l(\alpha) + \langle h_z, w \odot (X\theta - y) - z \rangle + \langle h_\alpha, \theta - \alpha \rangle \tag{35}$$

has a saddle point. Explicitly, there exist $(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*)$, not necessarily unique, for which

$$\mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) \leq \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) \leq \mathcal{L}_0(\theta, z, \alpha, h_z^*, h_\alpha^*) \tag{36}$$

holds for all $\theta, z, \alpha, h_z, h_\alpha$.

Proof The primitive problem is $\min_{\theta, z, \alpha} \sup_{h_z, h_\alpha} \mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha)$, represented by P^l . The dual problem is $\max_{h_z, h_\alpha} \inf_{\theta, z, \alpha} \mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha)$, represented by D^l . For $\mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha)$, since $f(\theta) + g(z) + l(\alpha)$ is a proper closed convex function, $w \odot (X\theta - y) - z = 0$ and $\theta - \alpha = 0$ is an affine function, and the existence points $(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*)$ satisfy the Karush–Kuhn–Tucker (KKT) condition, so according to the strong and weak duality and optimality conditions of the Lagrange multiplier method [24], the following conclusions can be obtained:

The primitive problem P^l is equal to the optimal value of the dual problem D^l , that is, $\text{val}(P^l) = \text{val}(D^l)$. The duality gap between the original problem and the dual problem is zero, which means that satisfies the strong max–min property, and P^l and D^l have the same optimal solution. Where $\text{val}(x)$ represents the value of x .

Any point $(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*)$ that satisfies the KKT condition in $\mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha)$ has

$$\inf_{\theta, z, \alpha} \mathcal{L}_0(\theta, z, \alpha, h_z^*, h_\alpha^*) \leq \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) \leq \sup_{h_z, h_\alpha} \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z, h_\alpha), \tag{37}$$

i.e.

$$\text{val}(D^l) \leq \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) \leq \text{val}(P^l). \tag{38}$$

When the duality gap between the primitive problem and the dual problem is zero, i.e. $\text{val}(P^l) = \text{val}(D^l)$, we can get

$$\mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) = \inf_{\theta, z, \alpha} \mathcal{L}_0(\theta, z, \alpha, h_z^*, h_\alpha^*) \leq \mathcal{L}_0(\theta, z, \alpha, h_z^*, h_\alpha^*), \quad \forall \theta, z, \alpha \in R^n. \tag{39}$$

The same reason can get

$$\begin{aligned} \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) &= \sup_{h_z, h_\alpha} \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z, h_\alpha) \\ &\geq \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z, h_\alpha), \quad \forall h_z, h_\alpha \in R^n. \end{aligned} \quad (40)$$

In summary

$$\begin{aligned} \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z, h_\alpha) &\leq \mathcal{L}_0(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*) \\ &\leq \mathcal{L}_0(\theta, z, \alpha, h_z^*, h_\alpha^*), \end{aligned} \quad (41)$$

that is, $\mathcal{L}_0(\theta, z, \alpha, h_z, h_\alpha)$ has a saddle point $(\theta^*, z^*, \alpha^*, h_z^*, h_\alpha^*)$, not necessarily unique. The standard Lagrangian function of Eq. (33) satisfies theorem 2 as evidence.

According to Theorem 1 and Theorem 2, the ADMM iteration satisfies the following conditions, and the convergence of proof Ref. [13] Appendix A:

Residual convergence. $r^k \rightarrow 0$ as $k \rightarrow \infty$, i.e., the iterates approach feasibility.

Objective convergence. $(\theta^k) + g(z^k) + l(\alpha^k) \rightarrow f(\theta^*) + g(z^*) + l(\alpha^*)$ as $k \rightarrow \infty$, i.e., the objective function of the iterates approaches the optimal value.

Dual variable convergence. $h_z^k \rightarrow h_z^*, h_\alpha^k \rightarrow h_\alpha^*$ as $k \rightarrow \infty$, where (h_z^*, h_α^*) is a dual optimal point.

We know that the convergence rate is another important concept, which reflects the convergence speed of an iterative algorithm. The authors of [25, 26] have shown that ADMM can achieve $O(1/k)$ global convergence, where k is the number of iterations, under a strong convexity assumption. Without this strong convexity assumption, the author of [27] gives the most general result of ADMM convergence speed. Their results only require that both objective-function terms are convex (not necessarily smooth). Since here $f(\theta)$, $g(z)$ and $l(\alpha)$ are both convex, using ADMM to solve SMLR problems can achieve $O(1/k)$ convergence.

6 Experiment

In this section, experiments will be conducted on several public face databases to demonstrate WHCSC performance.

6.1 Experimental settings

WHCSC is compared to existing related methods, including NMR, RSRC, Sparse Huber (SH), RCRC, IRGSC. For RSRC, parameter p defaults to 1, and τ takes the best of (0, 1). In SH, the parameter η defaults to 10. For WHCSC, the parameter p defaults to 1 and τ get the best between (0,1). The parameter p in IRGSC defaults to 1, and it should be noted that formula 21 in the IRGSC has errors and should be changed to

$$\begin{aligned} \min_s \sum_{i=1}^m \{s_i e_i^2 + \gamma s_i^2\} &= \min_s s + \frac{E^2}{2\gamma_2}, \\ \text{s.t. } s^T I &= 1, s_i \geq 0, i = 1 \sim m. \end{aligned} \quad (42)$$

where $E = [e_1^2, e_2^2, \dots, e_m^2]$, and the authors in [28] also have the same opinion. This article sets comparative experiments according to the original IRGSC article.

In Sect. 3, we described how to reduce intra-class changes (classification model I and classification model II) and increase the variation between classes (1 power exponent weights and 0.5 power exponent weights). In Sect. 6, we use WHCSC, RSRC, RCRC and its improved algorithm to experimentally test the two methods, other unspecified algorithms in accordance with the original essay method to help contrast. The 1 power exponent weights and 0.5 power exponent weights are tested for WHCSC, respectively, to prove the validity of the 3.2.1 and 3.2.2 theory, and the corresponding names are WHCSC_1 and WHCSC_0.5. RSRC tests the classification model I and classification model II, respectively, and the corresponding names are RSRC_1 and RSRC_2. The RCRC also tests the classification model I and classification model II, respectively, corresponding to RCRC_1 and RCRC_2. SH uses classification model II.

6.2 Face recognition without occlusion

The performance in WHCSC was first tested by illumination changes without occlusion. Datasets use ExYaleB database and PIE database.

1. FR with different samples size: This section tests the validity of WHCSC under changing the training sample size. The data set was randomly divided into two parts, one of which contained n images for each person for training and the other for testing, where $n = 10, 20, 30, 40, 50$. The already-divided data is saved to ensure that the different algorithm training sets and test sets are the same, and the average recognition rate of the 10 runs is counted. PIE database recognition rate as shown in Table 1, ExYaleB database recognition rate in Table 2. We can observe that WHCSC achieves the highest recognition rate in all other tests of ExYaleB database and PIE database except RSRC_1 and RSRC_2 at sample size 10. When the sample is larger than 30, WHCSC_1 is marginally higher than WHCSC_0.5. Second, the classification rates of RSRC_2, RCRC_2, SH and so on are higher than that of RSRC_1, RCRC_1, and NMR. In addition, RSRC is better than RCRC in most cases, reflecting the validity of its weight vector. Overall, the WHCSC proposed in this paper achieved the best results.

Table 1 PIE database recognition rates for different sample sizes of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, NMR (Unit: percentage)

Sample size	10	20	30	40	50
WHCSC_1	79	91	96.24	96.29	97.12
WHCSC_0.5	79.18	91.06	96.24	96.12	97.06
RSRC_1	86.23	89.17	90.07	92.24	93
RSRC_2	79.23	90.76	96.17	96.24	96.65
RCRC_1	83.92	87.62	90.11	92.1	93.14
RCRC_2	77.82	90.76	95.82	96.11	96.6
IRGSC	66.41	84.74	93.12	94.18	95.94
SH	78.64	91.01	96.17	96.17	97.06
NMR	77.12	89.34	92.63	93.06	94.87

Bold numbers in the table indicate the values that work best under this parameter

2. FR with different feature dimension: This section tests WHCSC performance under different feature dimension. For databases (ExYaleB and PIE), 20 samples per subject were randomly selected for training, the rest of the samples were used for testing. Saving the divided data to ensure that when the parameters are changed, the test data sets of different algorithms are the same, and the average recognition rate of 10 runs is counted. PCA is a recognized projection technique used to reduce the dimensions of the original face image [3]. From Tables 3 and 4, not all WHCSCs achieve the best recognition rate in the different dimensional character tests. All recognition rates for WHCSC_1 are better than WHCSC_0.5. The results of RCRC_2 and RSRC_2 in different feature dimensions are better than RCRC_1 and RSRC_1, respectively. In the test of different characteristic dimensions, not all algorithms reduce the recognition rate as the feature dimension decreases. For example, the recognition rate in the 200-dimension is mostly higher than the 150-dimensional and 250-dimensional. This is because after the PCA reduces the dimension, the feature tries to obtain a more meaningful low-dimensional representation, but in fact may lose the original dictionary information contained in the high-dimensional feature.

6.3 Face recognition with occlusion

One of the advantages of WHCSC is its robustness in terms of occlusion and noise damage. On the one hand, the parameter $W\eta$ is used to evaluate $g(z)$ meets l_2 loss or l_1 loss, thus reducing the influence of noise or outliers. On the other hand, classification model II and 1 power exponent weight make it easier to distinguish between different

Table 2 ExYaleB database recognition rates for different sample sizes of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, NMR (Unit: percentage)

Sample size	10	20	30	40	50
WHCSC_1	79	91.72	96.12	96.17	97.28
WHCSC_0.5	78.51	91.54	94.11	94.52	97.08
RSRC_1	85.94	87.3	90.73	91.67	94.94
RSRC_2	78.31	91.65	94.11	94.52	97.08
RCRC_1	83.92	87.85	90.73	94.53	94.94
RCRC_2	77.82	90.76	95.82	96.12	96.65
IRGSC	70.25	86.34	90.66	92.39	94.36
SH	77.48	90.8	94.19	94.63	96.5
NMR	76.25	90.34	92.5	93.89	94.02

Bold numbers in the table indicate the values that work best under this parameter

Table 3 PIE database recognition rates for different feature dimension of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, NMR (Unit: percentage)

Feature dimension	50	100	150	200	250	300
WHCSC_1	85.69	88.92	88.13	89.12	88.04	89.31
WHCSC_0.5	71.47	81.67	85	89.02	87.94	89.12
RSRC_1	64.41	81.04	82.11	84.57	83.42	83.45
RSRC_2	72.84	82.25	84.61	86.47	85.2	85.78
RCRC_1	65.2	82.25	86.27	87.06	88.3	89.41
RCRC_2	85.98	88.72	88.43	88.73	87.55	88.63
IRGSC	79.5	85.19	85.29	84.21	83.63	84.8
SH	85.88	88.72	88.33	88.92	87.74	89.22
NMR	68.82	76.11	82.14	82.1	83.5	84.11

Bold numbers in the table indicate the values that work best under this parameter

Table 4 ExYaleB database recognition rates for different feature dimension of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, NMR (Unit: percentage)

Feature dimension	50	100	150	200	250	300
WHCSC_1	87.85	89.54	90.63	92.93	91.47	91.54
WHCSC_0.5	87.84	89.48	90.81	92.86	91.29	91.17
RSRC_1	73.15	86.21	89.9	92.62	91.23	90.56
RSRC_2	85.97	88.87	90.93	92.38	91.83	91.17
RCRC_1	73.16	86.21	90.38	92.8	92.32	93.23
RCRC_2	88.33	89.29	90.75	92.14	91.29	90.87
IRGSC	80.77	81.98	84.64	86.03	84.95	85.61
SH	88.03	89.36	90.81	92.86	91.47	91.41
NMR	78.6	82.6	90.5	91.41	88.67	87.93

Bold numbers in the table indicate the values that work best under this parameter

Fig. 5 From left to right, Sample subset 1 through subset 5 sample images, respectively

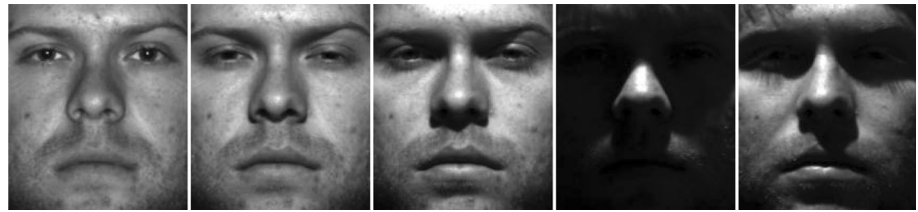


Fig. 6 Different percentage pixels damaged face images (from 0 to 70%)

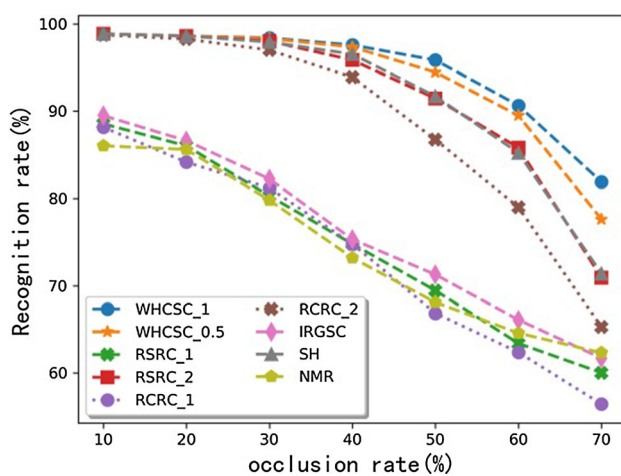


Fig. 7 Different pixel corrosion face recognition

categories of faces by lowering intra-class different and increasing inter-class difference. In this section, we will evaluate the robustness of WHCSC to different types of occlusions, such as Gaussian noise random pixel corruption, random block occlusion, masquerading, and so on. The WHCSCs are compared to existing related methods, including NMR, RSRC, Sparse Huber (SH), RCRC, IRGSC, and both the RSRC and RCRC will test both classification modes. The robustness of RSRC is achieved by repeatedly assigning weights to the training samples through a sigmoid function with variable parameters. The robustness of RCRC is achieved by sparse coding constrained coding coefficients. The robustness of SH is achieved by using l_2 loss and l_1 loss in combination with coding residuals. NMR is a recently proposed matrix-based regression classification method, which not only retains the structural information of face images but also has good robustness. The IRGSC achieves robustness by adaptive feature weights and distance-weighted learning. The

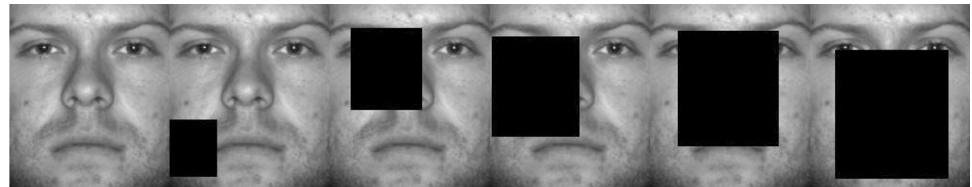
WHCSC robustness was tested by real complex occlusion experiments.

1. FR with pixel corrosion: This section uses the ExYaleB database, which has a total of 64 face images for each theme and can be divided into 5 subsets depending on the lighting conditions and face angle. Sample images of each subset are shown in Fig. 5, wherein subset 1 and subset 2 have good lighting conditions, subset 3 has medium lighting conditions, subset 4 has most poor lighting conditions, and subset 5 has poor lighting conditions. A total of 22 face images of our fixed-sampling sub-sets 1, 2, 3 and 5 were used for training, and the rest of the 4 subsets were used for testing. All images are cropped to 32×28 pixel size. For each test image, a fixed proportion of noise is added using random grayscale and random locations, i.e., Gaussian noise. The original image shown in Fig. 6 is a face image of 192×168 pixels with different pixel noise.

As can be observed in Fig. 7, the WHCSC test results are superior to other algorithms for pixel etches at different scales. Second, the recognition rate of the algorithm using classification model II is much higher than that of the algorithm using classification model I. The recognition rate of WHCSC_1 was 0.23%, 1.39%, 1.15% and 4.27% higher than that of WHCSC_0.5 when the signal to noise ratio was equal to 40%. In addition, the recognition rates of RSRC_1 and RSRC_2 are mostly higher than those of RCRC_1 and RCRC_2, respectively, and most of IRGSCs are better than RSRC_1. This indirectly verifies the validity of the IRGSC and RSRC algorithms. In summary, pixel-corrupted face recognition once again validates the robustness and validity of WHCSC for outliers. And on the other hand, it also validated the noise-based advantages of the classification model II and 1 power exponent weight.

2. FR with Block Occlusion: In this section, we design two block occlusion experiments. In the first

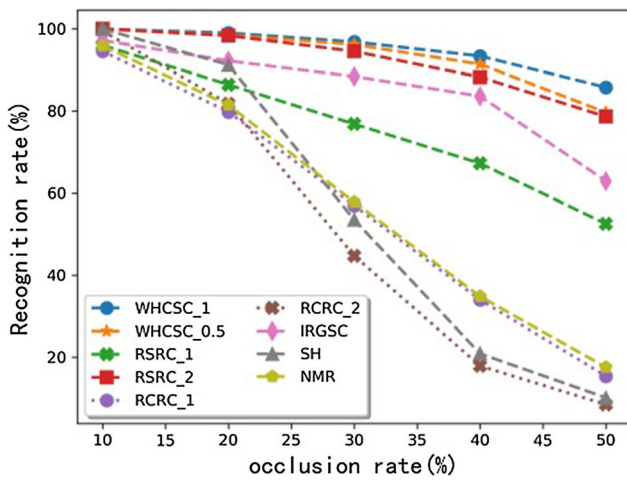
Fig. 8 **a** Face image of a black block occlusion, **b** face image of a white block occlusion



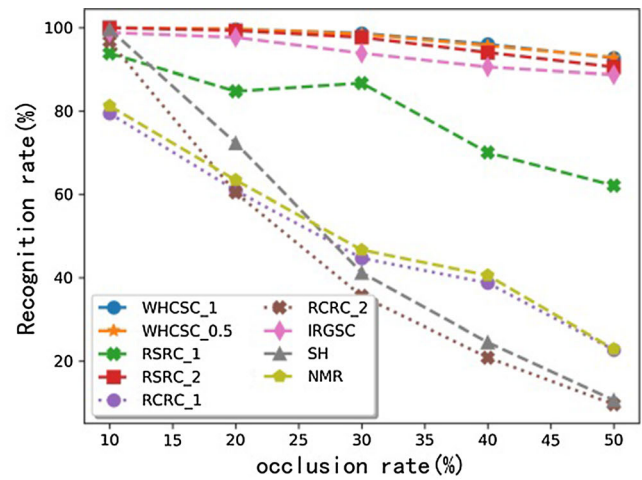
(a)



(b)



(a)



(b)

Fig. 9 **a** Face recognition rate of black block occlusion, and **b** face recognition rate of white block occlusion

Fig. 10 The face image with Lina block occlusion



experiment, we replaced 10–50% of each test image with white or black blocks. Half of the face images of the fixed subset 1, 2 and 3 were acquired for training and the rest of the 3 subsets were used for testing. The position of the occlusion box is random. Figure 8 shows a partially occluded facial image of the ExYaleB database with different block blocking ratios. Figure 9 shows the recognition rates of RSRC, RCRC, IRGSC, SH, NMR and WHCSC in different block occlusions. We can observe that WHCSC has the obvious advantage of having the highest recognition rate at different

occlusion percentages. At occlusion percentages above 20%, the RCRC, SH, and NMR discrimination rates dropped significantly. The recognition rate of WHCSC was 86.76% when the shielding ratio of black block reached 50%, which was 7.12% higher than that of RSRC_2 and 22.85% higher than that of IRGSC. However, RCRC, SH, and NMR had failed at this time. Meanwhile, WHCSC_1 is 6.12% higher than WHCSC_0.5, and RSRC_2 is 26.16% higher than RSRC_1. WHCSC_1 has a 92.72% recognition rate when the white block occlusion ratio reaches 50%,

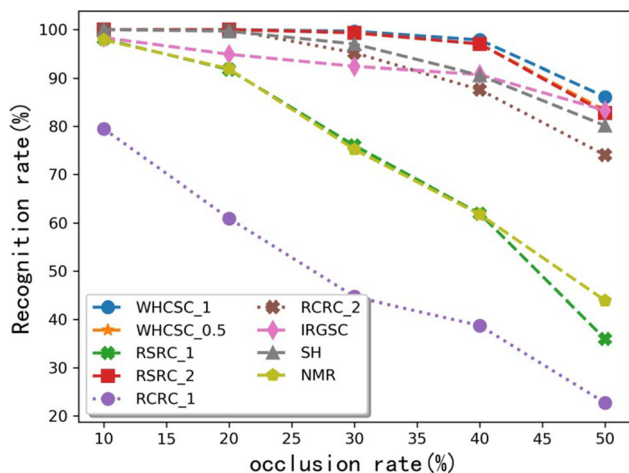


Fig. 11 The face recognition rate with Lina block occlusion

Table 5 Recognition rate (unit: percentage) of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, NMR in sunglasses camouflage and scarf camouflage

	WHCSC_1	WHCSC_0.5	RSRC_1	RSRC_2
Sunglasses	93.67	92	39.33	84.67
Scarf	86	82.33	38.33	51.67

	RCRC_1	RCRC_2	IRGSC	SH	NMR
Sunglasses	23.83	45.83	77.67	51.5	23.67
Scarf	26.5	15.64	62.83	7.33	27.33

which is 2.16% higher than RSRC_2 and 3.98% higher than IRGSC. Meanwhile, RSRC_2 is 28.47% more than RSRC_1. However, RCRC, SH, NMR still failed. WHCSC_1 has the best occlusion ratio except 0.16% lower than WHCSC_0.5 at 50% occlusion percentage.

In the second experiment, the classic Lena diagram was used as the occlusion element to replace 10–50% pixel of each test image. Figure 10 shows the test image samples. We can see that the pixels in the occlusion area are close to the original pixels relative to the first two experiments. Figure 11 shows the recognition rates of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH and NMR at 10–50% block occlusion. It can be observed that the overall recognition rate is increasing, and WHCSC still maintains the highest recognition rate. Surprisingly, RCRC_2 and SH showed a good recognition rate. On the one hand, it is easier to train a linear combination of images because of the occlusion area close to the original pixel. On the other hand, we further prove the good effect of WHCSC weight on image local optimization.

3. Real camouflage face recognition: The experiment in this section uses AR database, using the first three of

each face subset 1 and subset 2, a total of six as a training image. Six pieces of camouflage images in Subset 1 and Subset 2 and 6 pieces of scarf camouflage images were taken as test images, respectively. The image is adjusted to 33×24 pixels. Table 5 shows the test results of several classifiers, WHCSC shows better results than RSRC, RCRC, IRGSC, SH, NMR. The performance of RCRC is unstable because of scarf camouflage masks effective pixels of more people, making RCRC vulnerable to interference when the image information is limited. IRGSC performs well and further reflects the effect of weight coefficient on local image optimization.

6.4 Image reconstruction

Reconstructed block occlusion and real camouflage fitted images, and observe the reconstructed image of each algorithm. In this section, the training set uses a frontal, non-occluded image, and the test set uses an occlusion image. In algorithms with weights coefficients such as WHCSC, RSRC, IRGSC, the reconstructed image is represented as $w \odot X\theta$, and the corresponding test set image is represented as $w \odot y$. In algorithms without weights coefficients such as RCRC, SH, NMR, the reconstructed image is represented as $X\theta$, and the corresponding test set image is represented as y . The noise in the pixel-corroded image is randomly distributed, and the reconstructed image is not easy to observe, so no experiment is set.

Figure 12 is an image reconstruction of block occlusion. Looking at (f) and (g) of Fig. 12, because of a black occlusion block in the test set, these algorithms without weights coefficients in the reconstructed image, such as RCRC, SH and NMR, cannot generate an area similar to a black occlusion block. The reconstructed image of RSRC_1 has begun to corrode other normal image areas in a large amount when the black block has not been completely fitted. WHCSC_1, WHCSC_0.5, RSRC_2 and IRGSC can all fit the black occlusion block well. A closer look reveals that the forehead of the test image has subtle color differences due to different illumination angles. When the black occlusion block is completely fitted, RSRC_2 has obvious noise corrosion in the forehead area of the face. WHCSC_0.5 and IRGSC have slight noise corrosion, and WHCSC_1 is almost none. The (h) of Fig. 12 shows that as the parameter τ decreases, the residual threshold η is smaller and the weights constraint is stronger. When τ is equal to 0.9, 0.8, 0.7 and 0.6, the reconstructed image does not completely fit the black block area; and when τ is 0.5, the η that is too small causes the weight to over-constrain the residual, thus causing corrosion of the pixels outside the black block area. When τ is

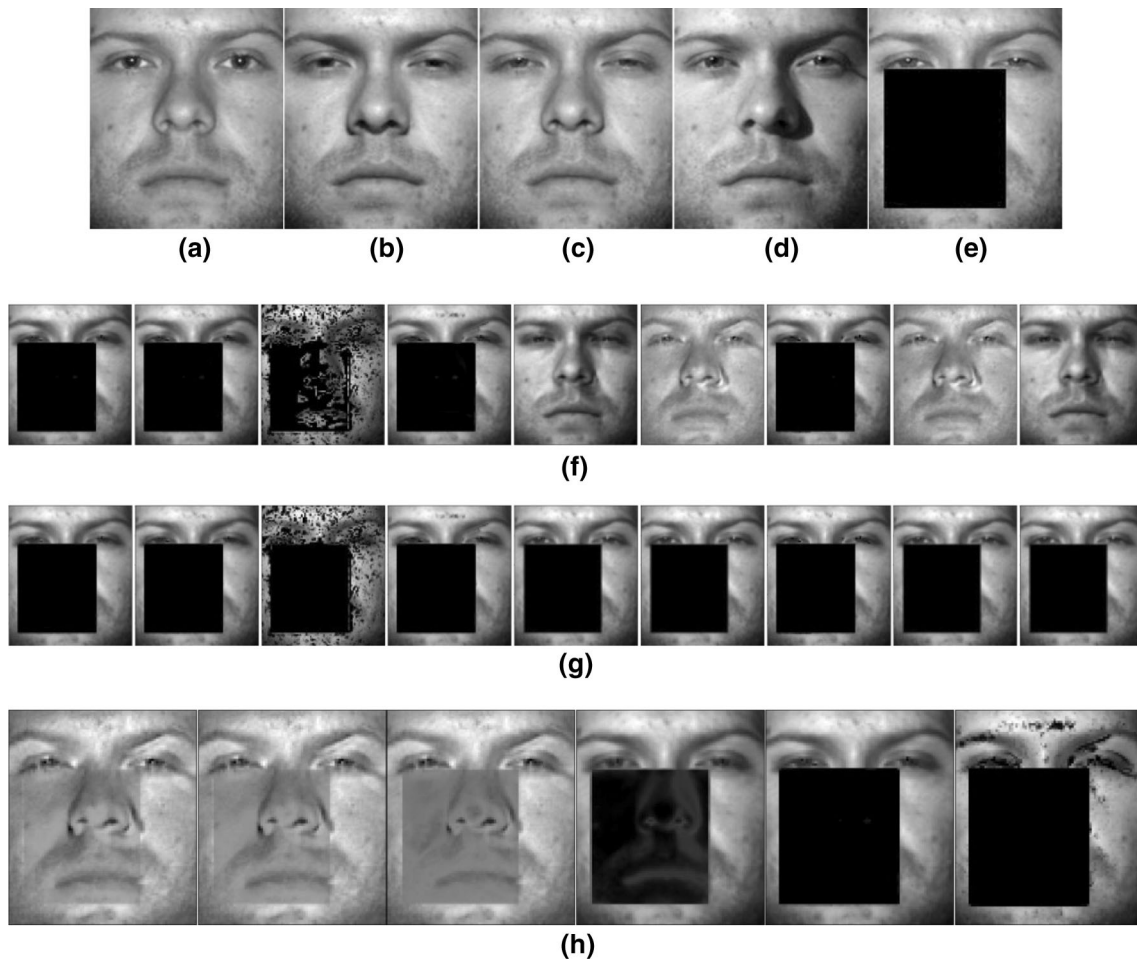


Fig. 12 Reconstructed image of a 40% black block occluded face. From **a–d** are training sets. **e** Test sets. In **f**, reconstructed images of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, and NMR are shown from left to right. **g** shows a

comparison of test sets of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH, and NMR from left to right. **h** Reconstructed image of WHCSC_1 when the parameter τ is equal to 0.9, 0.8, 0.7, 0.6, 0.58, and 0.5, respectively

0.58, the reconstructed image completely fits the black block area, and almost no other pixel points are corroded. In addition, $\tau = 0.58$ indicates that 42% of the pixels are considered to have larger residuals, slightly larger than the test set by 40%. Because the real image itself has noise generated by other factors, it is in line with theory and practice to reconstruct the image to obtain the optimal performance at $\tau = 0.58$ summary, in a complex noisy environment, the parameter q can make the weight coefficient smoother, and the value of the parameter τ can be easily determined by the actual number of residuals, which continues to show the superiority of the WHCSC.

Figure 13 is the image reconstruction of the sunglasses camouflage. Looking at (h) and (i) of Fig. 13, since there is a sunglasses camouflage in the test sets, the algorithms with no weights cannot generate a region similar to the sunglasses in the reconstructed image, such as RCRC, SH and NMR. The reconstructed image of RSRC_1 has only a faint sunglasses frame, and the entire image is cluttered with

noise. This indicates that in the noisy environment, the classification model I does not distinguish between noise and real images very well. Although IRGSC fits the sunglasses camouflage, its weights coefficient is not accurate enough for the boundary of the character’s outline and expression.

Compared with the corresponding test set comparison chart, the reconstructed images of WHCSC_1, WHCSC_0.5 and RSRC_2 have almost no difference, and they can reconstruct the characteristics of the test image very well. On the one hand, the reconstructed image fits the shape and gloss of the sunglasses. For the test set, the white area on the sunglasses belongs to the image point with large residuals, and its pixel value approaches zero under the weight coefficient, that is, it is black in the gray image. On the other hand, the reconstructed image weakens the facial expression influences brought about by (b) in Fig. 13.

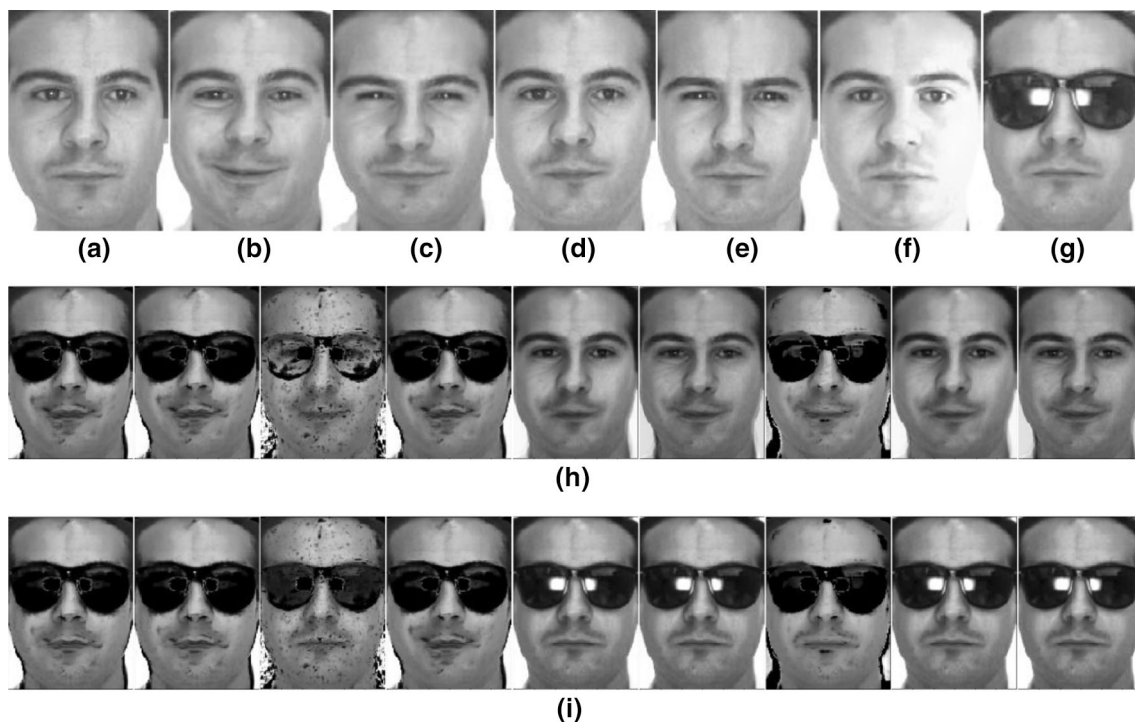


Fig. 13 Face reconstruction image with sunglasses camouflage. From **a–f** is a training sets for sunglasses camouflage. **g** Test sets for sunglasses camouflage. Reconstructed images of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH,

and NMR are shown from left to right in **h**. In **i**, a comparison images of test sets of WHCSC_1, WHCSC_0.5, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH and NMR are shown from left to right

Table 6 Run-time tests such as WHCSC_1, RSRC_1, RSRC_2, RCRC_1, RCRC_2, IRGSC, SH. Unit: second

	Gaussian noise		Sunglasses camouflage	
	Correct rate	Running time	Correct rate	Running time
WHCSC_1	81.89	572	93.17	657
RSRC_1	60.03	942	39.33	430
RSRC_2	70.93	700	84.67	894
RCRC_1	56.43	1068	23.83	449
RCRC_2	65.28	739	45.83	619
IRGSC	61.7	2307	77.67	1219
SH	71.39	541	51.5	261

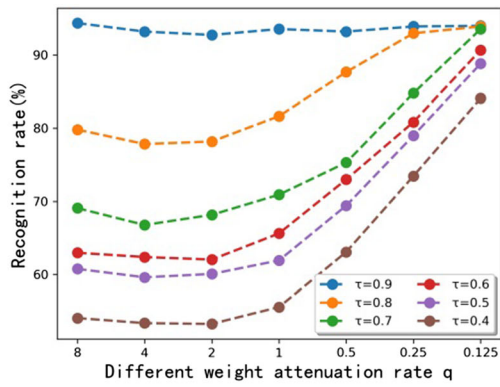
6.5 Running time

Running time is one of the important reference indexes for judging classifier. Five robust classifiers, such as WHCSC, RSRC, RCRC, IRGSC, and SH, were run on the same computer with noise and real disguise. Algorithms involving l_1 norm minimization are all solved by ADMM. Table 6 lists the average run times for the 10 runs of the 5 classifiers. The IRGSC takes the longest time due to the additional computation of adaptive feature weights and adaptive distance weights, and its recognition rate is medium and stable; SH takes the least time, but the recognition rate is low; WHCSC_1 is less and the recognition rate is the highest and stable. Influenced by the

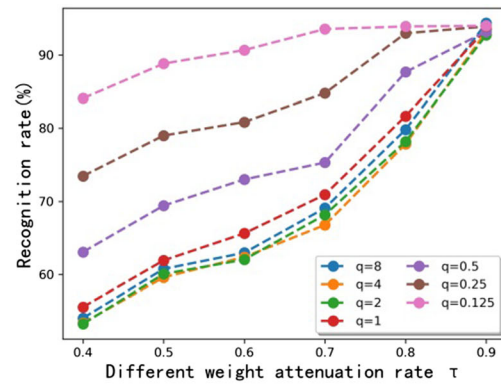
experimental samples, the computational cost of classification model II has no obvious advantage. In summary, WHCSC sacrificed a small amount of computational cost and achieved the highest recognition rate.

6.6 The impact of parameters on the recognition rate

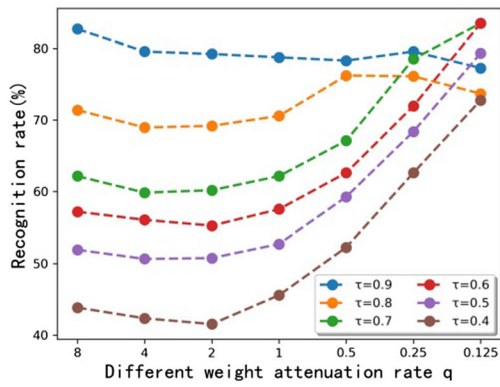
Parameter changes are another important reference indicator for judging classifiers. WHCSC has two important parameters, such as τ and q mentioned earlier. The position of the threshold residual in the residual sequence Ψ is determined by $k = \lfloor \tau m \rfloor$; and the penalty rate of the weight is affected by q .



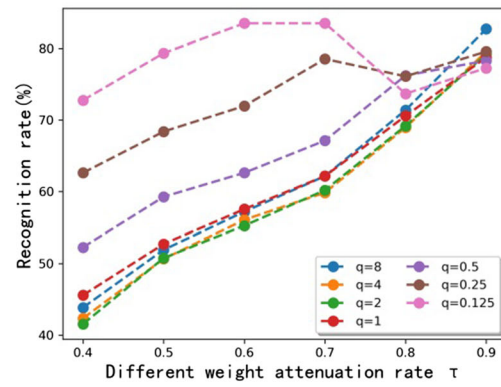
(a)



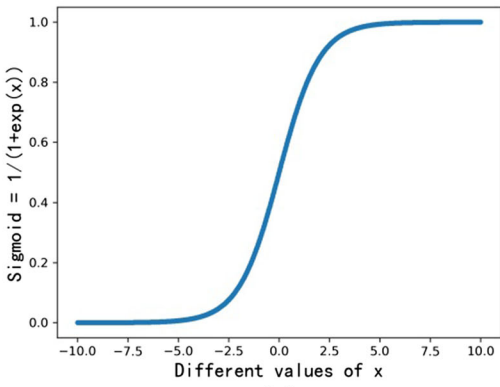
(b)



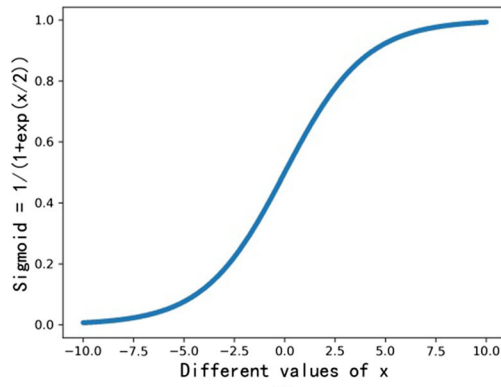
(c)



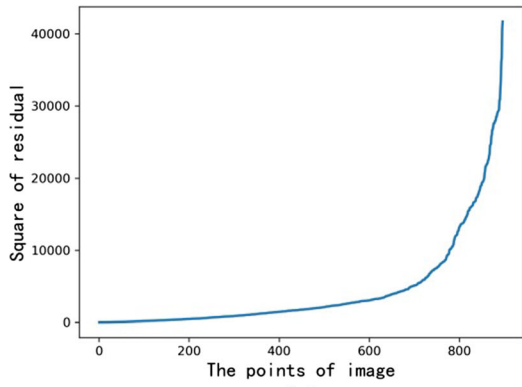
(d)



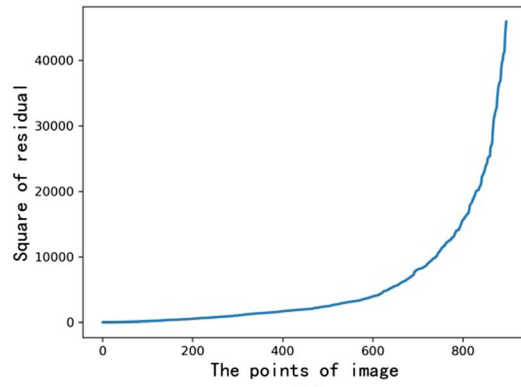
(e)



(f)



(g)



(h)

◀ **Fig. 14** **a** Recognition rate of the different parameters q in the case of 50% Gaussian noise. **b** Recognition rate of different parameters τ in the case of 50% Gaussian noise. **c** Recognition rate of the different parameters q in the case of 70% Gaussian noise. **d** Recognition rate of different parameters τ in the case of 70% Gaussian noise. **e** Schematic diagram of $\text{sigmoid} = \frac{1}{1+\exp(x)}$. **f** Schematic diagram of $\text{sigmoid} = \frac{1}{1+\exp(x/2)}$. **g** Residual distribution map of 50% Gaussian noise. **h** Residual distribution map of 70% Gaussian noise

(a) and (c) of Fig. 14 show the fixed parameter τ , and the change in the recognition rate when the parameter q is changed. As q decreases, the recognition rate shows an upward trend overall. (b) and (d) of Fig. 14 show the change of the recognition rate when the parameter q is fixed and the parameter τ is changed. As τ increases, the recognition rate shows an upward trend overall. (e) and (f) of Fig. 14 show some of the characteristics of the sigmoid function. When the parameter x changes to $\frac{x}{2}$, the sigmoid function image is smoother. Therefore, when the parameter q becomes small, the degree of weights penalty can be reduced, so that the value of the weights change trend is smoother in the same iteration. (g) and (h) of Fig. 14 show residual distribution maps of face images of 896 size. It can be observed that only a small number of face images have large residuals. Therefore, when the parameter τ is increased, more image points can be obtained to obtain higher weights. In combination with the complex environment of the face image, when the noise is enhanced, the image points with large residuals will also increase, and the value of τ should be lowered. Conversely, the value of τ can be increased. The parameter q is usually small.

7 Conclusion

In this paper, we propose a newly weighted Huber constrained sparse coding, and propose an effective optimization method to enhance the effect of weights. The benefits of WHCSC are reflected in the robustness and effectiveness of the occluded complex environment. On the one hand, the weight constraint can effectively find the noise pixels in the query sample and reduce the weight of the noise pixels at the time of regression, which achieves the local optimization. On the other hand, we use Huber's estimation to choose different fidelity terms (l_1 or l_2 norm) to further accurately return the query samples. Secondly, the use of classification mode II can avoid the interference caused by other types of images when the current category is regressed. Finally, increasing the variability between classes through 1 power exponent weight makes it easier to classify. WHCSC is suitable for complex changes of PCA, illumination, corrosion, and occlusion. Experiments show

that WHCSC is superior to IRGSC, RSRC, SRC, NMR and so on, and it is smoother and more accurate for noise processing. Its high robustness and strong effectiveness are the ideal choices in face recognition applications.

Acknowledgement This paper is supported by the following foundations or programs, including Chongqing Innovative Project of Overseas Study (No. cx2018120), National Social Science Foundation of China (No. 17XFX013). The authors would like to thank the anonymous referees for their valuable comments and suggestions.

Compliance with ethical standards

Conflict of interest We declare that we have no financial and personal relationships with other people or organizations that can inappropriately influence our work, there is no professional or other personal interest of any nature or kind in any product, service and/or company that could be construed as influencing the position presented in, or the review of, the manuscript entitled, "Weighted Huber Constrained Sparse Face Recognition".

References

- Liu L, Xiong C, Zhang H, Niu Z, Wang M, Yan S (2016) Deep aging face verification with large gaps. *IEEE Trans Multimed* 18(1):64–75
- Kan M, Shan S, Zhang H, Lao S, Chen X (2015) Multi-view discriminant analysis. *IEEE Trans Pattern Anal Mach Intell* 38(1):188–194
- Jiang X (2009) Asymmetric principal component and discriminant analyses for pattern classification. *IEEE Trans Pattern Anal Mach Intell* 31(5):931–937
- Tolba AS, El-Baz AH, El-Harby AA (2006) Face recognition: a literature review. *Int J Signal Process* 2(1):88–103
- Xu Y, Li Z, Zhang B, Yang J, You J (2017) Sample diversity, representation effectiveness and robust dictionary learning for face recognition. *Inf Sci* 375:171–182
- Naseem I, Togneri R, Bennamoun M (2010) Linear regression for face recognition. *IEEE Trans Pattern Anal Mach Intell* 32(11):2106–2112
- Wright J, Yang AY, Ganesh A, Sastry SS, Ma Y (2009) Robust face recognition via sparse representation. *IEEE Trans Pattern Anal Mach Intell* 31(2):210–227
- Zhang L, Yang M, Feng X (2011) Sparse representation or collaborative representation: which helps face recognition? In: 2011 international conference on computer vision, Barcelona, pp 471–478
- Jian Y, Lei L, Qian J, Ying T, Zhang F, Yong X (2017) Nuclear norm based matrix regression with applications to face recognition with occlusion and illumination changes. *IEEE Trans Pattern Anal Mach Intell* 39(1):156–171
- Zhong D, Xie Z, Li Y, Han J (2015) Loose $L_{1/2}$ regularised sparse representation for face recognition. *Comput Vis IET* 9(2):251–258
- Zheng J, Yang P, Chen S, Shen G, Wang W (2017) Iterative re-constrained group sparse face recognition with adaptive weights learning. *IEEE Trans Image Process* 26(5):2408–2423
- Lin G, Yang M, Yang J, Shen L, Xie W (2018) Robust, discriminative and comprehensive dictionary learning for face recognition. *Pattern Recogn* 81:341–356
- Boyd S, Parikh N, Chu E, Peleato B, Eckstein J (2010) Distributed optimization and statistical learning via the alternating

- direction method of multipliers. *Found Trends Mach Learn* 3(1):1–122
14. Zhang H, Wang S, Xu X, Chow TWS, Wu QMJ (2018) Tree2Vector: learning a vectorial representation for tree-structured data. *IEEE Trans Neural Netw Learn Syst* 29(11):1–15
 15. Zhang H, Wang S, Zhao M, Xu X, Ye Y (2018) Locality reconstruction models for book representation. *IEEE Trans Knowl Data Eng* 30(10), pp Locality reconstruction models for book representation, 2018
 16. Olshausen BA, Field DJ (2004) Sparse coding of sensory inputs. *Curr Opin Neurobiol* 14(4):481–487
 17. Földiák P, Young MP (1998) Sparse coding in the primate cortex. In: Michael AA (ed) *The handbook of brain theory and neural networks*. MIT Press, Cambridge, pp 895–898
 18. Liu Y, Li X, Liu C, Liu H (2016) Structure-constrained low-rank and partial sparse representation with sample selection for image classification. *Pattern Recogn* 59:5–13
 19. Liu Y, Wu F, Zhang Z, Zhuang Y, Yan S (2010) Sparse representation using nonnegative curds and whey. In: 2010 IEEE conference on computer vision and pattern recognition, San Francisco, CA, pp 3578–3585
 20. Gao S, Tsang WH, Chia LT, Zhao P (2010) Local features are not lonely—Laplacian sparse coding for image classification. In: 2010 IEEE computer society conference on computer vision and pattern recognition, San Francisco, CA, pp 3555–3561
 21. Wang J, Yang J, Yu K, Lv F, Huang T, Gong Y (2010) Locality-constrained linear coding for image classification. In: 2010 IEEE computer society conference on computer vision and pattern recognition, San Francisco, CA, pp 3360–3367
 22. Ramírez I, Lecumberry F, Sapiro G (2009) Universal priors for sparse modeling. In: 2009 3rd IEEE international workshop on computational advances in multi-sensor adaptive processing (CAMSAP), Aruba, Dutch Antilles, pp 197–200
 23. Yang M, Zhang L, Yang J, Zhang D (2011) Robust sparse coding for face recognition. In: 2011 IEEE conference on computer vision and pattern recognition, pp. 625–632
 24. Boyd S, Vandenberghe L, Foybusovich L (2006) Convex optimization. *IEEE Trans Autom Control* 51(11):233–237
 25. Deng W, Yin W (2015) On the global and linear convergence of the generalized alternating direction method of multipliers. *J Sci Comput* 66(3):889–916
 26. Goldstein T, Donoghue B, Setzer S (2014) Fast alternating direction optimization methods. *SIAM J Imaging Sci* 7(3):225–231
 27. He B, Yuan X (2012) On non-ergodic convergence rate of Douglas–Rachford alternating direction method of multipliers. *Numer Math* 130(3):567–577
 28. Li H, Hu H, Yip C (2017) Comments on “Iterative re-constrained group sparse face recognition with adaptive weights learning”. *IEEE Trans Image Process* 26(11):5475–5476

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.