

Jointly Sparse Locality Regression for Image Feature Extraction

Dongmei Mo , Zhihui Lai , Xizhao Wang , *Fellow, IEEE*, and Waikeung Wong 

Abstract—This paper proposes a novel method called Jointly Sparse Locality Regression (JSLR) for feature extraction and selection. JSLR utilizes joint $L_{2,1}$ -norm minimization on regularization term, and also introduces the locality to characterize the local geometric structure of the data. There are three main contributions in JSLR for face recognition. Firstly, it eliminates the drawback in ridge regression and Linear Discriminant Analysis (LDA) that when the number of the classes is too small, not enough projections can be obtained for feature extraction. Secondly, by using the local geometric structure as the regularization term, JSLR is able to preserve local information and find an embedding subspace which can detect the most essential data manifold structure. Moreover, since the $L_{2,1}$ -norm based loss function is robust to outliers in data points, JSLR provides the joint sparsity for robust feature selection. The theoretical connections of the proposed method and the previous regression methods are explored and the convergence of the proposed algorithm is also proved. Experimental evaluation on several well-known data sets shows the merits of the proposed method on feature selection and classification.

Index Terms—Regression, face recognition, feature extraction, local structure, joint sparsity.

I. INTRODUCTION

SINCE the data used in computer vision or pattern recognition is very high dimensional, it is of great importance to select the key features from large quantities of variables. Besides, the redundancy of the data would affect the performance of some algorithms in practical applications [1], and thus most of the algorithms cannot obtain a good performance in high-dimensional case [2]. Therefore, feature extraction and selection are of great importance in processing the high-dimensional data set [3].

Manuscript received May 4, 2019; revised September 9, 2019 and December 11, 2019; accepted December 13, 2019. This work was supported in part by The Hong Kong Polytechnic University (Project Code: RHR1) and in part by General Research Fund of the Research Grants Council of Hong Kong (Project Code: 15202217). The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jingdong Wang. (*Corresponding author: Waikeung Wong*.)

D. Mo and W. Wong are with the Institute of Textiles and Clothing, The Hong Kong Polytechnic University, Hong Kong SAR of China (e-mail: dongmei.mo@connect.polyu.hk; calvin.wong@polyu.edu.hk).

Z. Lai is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China, and also with the Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen 518060, China (e-mail: lai_zhi_hui@163.com).

X. Wang is with the College of Computer Science and Software Engineering and Guangdong Key Lab of Intelligent Information Processing, Shenzhen University, Shenzhen 518060, China (e-mail: xizhaowang@ieee.org).

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2019.2961508

Up to now, one of the classical methods is Principle Component Analysis (PCA) [4], which is a simple and effective unsupervised method as it solves the eigen decomposition problem to obtain the optimal vectors for dimensionality reduction. Linear Discriminant Analysis (LDA) [5] is a representative supervised method in feature extraction and dimensionality reduction, which uses the label information to improve the performance in classification. By maximizing the ratio of the between-class scatter to the within-class scatter of the training dataset, LDA can obtain an optimal set of discriminative vectors [6]. However, a drawback of LDA is that it is unsuitable for small sample size problem in face recognition. An effective model called PCA+LDA [7], which joints the two major techniques to obtain the discriminant vectors [8], has been proposed to deal with the problem. The other two methods called Particle Swarm Optimization (PSO) [9] and Backtracking Search Algorithm (BSA) [10] can also significantly reduce the number of features so as to reduce the computational complexity and at the same time guarantee the same level of performance.

However, PCA and LDA cannot provide the sparse projections for feature extraction since the learned projections are the linear combination of the data [11]. Recently, sparse regression showed the outstanding performance in feature selection and extraction [12]–[15]. By adding sparsity penalty for feature selection, the accuracy and robustness of these methods might be improved. Thus many studies focused on the sparse learning for variable selection. Zou *et al.* proposed an effective model called Sparse Principle Component Analysis (SPCA) [4] to generate modified principle components with sparse loadings by using the lasso or elastic net constraint [16], [17]. Some other sparse PCA algorithms, such as the SCoTLASS algorithm [18], the DSPCA algorithm [19] were proposed. All of these methods focused on sparse learning without using the class label information. Besides, Feng *et al.* proposed the unsupervised learning method based on maximum information and minimum redundancy (MIMR) [20] for hyperspectral image analysis, and Li *et al.* proposed another unsupervised feature selection method by nonnegative spectral analysis and redundancy control [21].

Some other researchers developed the supervised methods using the label information to perform sparse learning for feature extraction and selection. One of the effective methods is Sparse Discriminant Analysis (SDA) [22], which extends linear discriminant analysis to sparse case by imposing the sparsity constraint. Moreover, to overcome the data piling problem of LDA in the high dimensional and low sample size (HDLSS) case, Qiao *et al.* proposed sparse LDA to obtain sparse linear

discriminant vectors by taking the relationship between Fisher's LDA and a generalized eigenvalue problem into consideration [23]. Besides, some semi-supervised methods were also proposed. A semi-supervised method, which used partially labeled data samples, was designed to achieve batch feature selection [24]. Another method called Hessian sparse feature selection based on $L_{2,1/2}$ -matrix norm (HFSL) was proposed for semi-supervised sparse feature selection [25]. For the multimodal case, Ding *et al.* proposed a method using multimodal information to jointly learn face representation [26]. Furthermore, sparse regularization learning were also used in classification designation for different pattern recognition tasks [27]–[29].

It is a well-known fact that not all data are distributed on a linear subspace. They may lie on the nonlinear low-dimensional manifold embedding on the high-dimensional ambient space. Therefore, many manifold learning algorithms were proposed. The representative methods include Neighborhood Preserving Embedding (NPE) [30], Isometric Projection (IsoP) [31] and Locality Preserving Projection(LPP) [32], [33], etc. These algorithms aimed to preserve the local geometric structure of the data manifold. By introducing the locality for sparse subspace learning, Cai *et al.* also proposed a new method called Unified Sparse Subspace Learning (USSL) [34]. USSL utilized the elastic net for regression to simultaneously select the most important variables and take the local geometric structure into consideration. Besides, by combining the global pairwise sample similarity with local geometric structure, a new method called GLSPFS [35] was proposed by Liu *et al.* for feature selection.

In recent years, a great deal of attention has been paid to the regression methods with different norms for image recognition, feature extraction and variable selection [36]. For example, nuclear norm regression methods were proposed in [37], [38] for face recognition. The L_1 -norm based sparse regularized learning methods [39]–[41] have been used for face reconstruction and recognition [42]. A feature selection algorithm framework called Feature-weighting as Regularized Energy-based Learning (FREL) was proposed by Li *et al.* [43]. Based on least square regularization, Yang *et al.* [44] proposed the discriminative projection method. And the traditional RDA was further developed as Parameterless Reconstructive Discriminant Analysis (PRDA) [45] for feature extraction. In [46], the L_1 -norm minimization was employed to design a specific loss function, by which the abundant user tagged Web images are treated as noisy samples and will not be emphasized so as to perform robust semantic video indexing. Other methods, such as [47]–[53] were also proposed to deal with different feature selection problems. The methods in terms of jointly sparse subspace learning attracted great attention in the field of feature selection. Since the $L_{2,1}$ -norm based regression loss function is robust to outlier in data set, it can improve the robustness in learning steps. Therefore, some algorithms with joint $L_{2,1}$ -norm regularization were proposed to guarantee the joint sparsity for feature extraction. The model called Robust Feature Selection (RFS) [54] via joint $L_{2,1}$ -norms minimization showed the good performance for feature selection with joint sparsity. Yang *et al.* proposed another model called Unsupervised Discriminative Feature Selection (UDFS) [55] for sparse subspace learning. Experimental

results showed that UDFS outperforms the existing unsupervised feature extraction methods and its main advantage is that UDFS not only uses discriminative information but also uses local structure of datas distribution for feature selection [56]. The $L_{2,1}$ -norm regularization is also used in [57] to discover the common features shared across all the clustering tasks so as to obtain a discriminative low dimensional space for clustering. Except for the jointly sparse feature selection, the $L_{2,1}$ -norm was also widely used to deal with the joint-sparse recovery problems [58], [59] in computer vision.

Although a lot of methods have been developed to improve the performance of regression methods, there still exist some problems to be solved. For example, when the number of the class is too small, not enough projections can be obtained by the classical regression methods and/or their extensions to achieve higher classification accuracy. Also, most existing regression method do not simultaneously consider the geometric structure of the data as well as the sparsity of the projection matrix. In this paper, we propose a novel model called Jointly Sparse Locality Regression (JSR) for feature extraction and selection. JSR can not only avoid the limitation in the existing regression methods but also guarantee the sparsity by using $L_{2,1}$ -norm regularization on the projection matrix. What is more, JSR incorporates the local structure of the data in regression form, by which the optimization problem can be easily optimized so as to obtain better performance of feature extraction with less computational time.

The main contributions of this paper are described as below:

- 1) The number of the projections in LDA-based methods or regression-based methods is limited by the rank of the so-called between-class scatter matrix or the number of the classes. The proposed method can break out the limitation to obtain more projections for feature extraction by designing a novel regression model.
- 2) Theoretical connections between the proposed method and the previous regression methods are discovered. Moreover, the convergence of the proposed algorithm is also proved.
- 3) The experimental results of the proposed model with or without $L_{2,1}$ -norm regularization indicate that adding $L_{2,1}$ -norm penalty on the projection matrix can obtain joint sparsity for feature extraction so as to achieve high recognition rate.

The rest of this paper is organized as follows: In Section II, we discuss the related works and the extension based on ridge regression will be shown in Section III. In Section IV, we propose our objective function and the local optimal solution. Section V focuses on theoretical analysis (the convergence and the computational complexity). The proposed model will be evaluated by several well-known databases in Section VI. In Section VII, we draw a conclusion for this paper.

II. RELATED WORKS

In this section, the notations used in this paper will be briefly described and the related works will be reviewed.

192 *A. Notations*

193 Scalars are denoted as lowercase or uppercase italic letters,
 194 i.e. $i, j, d, p, n, c, \alpha_1, \alpha_2$ etc. while vectors are represented as
 195 bold lowercase italic letters, i.e. \mathbf{x}, \mathbf{y} , etc. Matrices are defined
 196 as bold uppercase italic letters, i.e. $\mathbf{A}, \mathbf{B}, \mathbf{X}, \mathbf{Y}, \mathbf{W}$ etc.

197 Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in R^{d \times n}$ then \mathbf{X} denotes a $d \times n$
 198 matrix as the original data set, where n is the number of total
 199 training samples and d denotes the features dimension for each
 200 sample. Let $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_c] \in R^{n \times c}$ to be a $n \times c$ matrix
 201 as the label matrix falling into c classes.

202 *B. Regressions*

203 Ridge Regression [60] is a regularized least square method
 204 for multivariate learning. It aims to solve the multicollinearity
 205 problem of covariates in samples.

206 The optimization problem of the simplest regression is

$$\mathbf{P}^0 = \arg \min_{\mathbf{P}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2 \quad (1)$$

207 where \mathbf{X} denotes the training set of n training data. The ma-
 208 trix $\mathbf{P} \in R^{d \times c}$ aims to lead the linear dependency between the
 209 training data and the corresponding labels. By setting the deriva-
 210 tives of (1) with respect to \mathbf{P} equaling to 0, we have the optimal
 211 solution

$$\mathbf{P}^0 = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \quad (2)$$

212 However this optimal solution is only suitable for the case when
 213 $\mathbf{X} \mathbf{X}^T$ is a full-rank matrix. Because of the small-sample size
 214 problem, the matrix $\mathbf{X} \mathbf{X}^T$ may be not a full-rank one. Therefore,
 215 to solve the singular problem in computing the inverse of
 216 $\mathbf{X} \mathbf{X}^T$ the L_2 -norm regularized term was added to (1), and then
 217 we have the classical ridge regression optimization problem:

$$\mathbf{P}^* = \arg \min_{\mathbf{P}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2 + \alpha \|\mathbf{P}\|_F^2 \quad (3)$$

218 By setting the derivatives of (3) with respect to \mathbf{P} equaling to 0,
 219 we have the optimal solution for (3) as

$$\mathbf{P}^* = (\mathbf{X} \mathbf{X}^T + \alpha \mathbf{I})^{-1} \mathbf{X} \mathbf{Y} \quad (4)$$

220 For further analysis in the following sections, we need to rep-
 221 resent the optimal solution of (3). Based on the SVD of $\mathbf{X} =$
 222 $\mathbf{U} \mathbf{D} \mathbf{V}^T$, the optimal solution can be represented as

$$\mathbf{P}^* = \mathbf{U} \frac{\mathbf{D}}{\mathbf{D}^2 + \alpha \mathbf{I}} \mathbf{V}^T \mathbf{Y} \quad (5)$$

223 From (2) and (5), we can know that the optimal projection matrix
 224 \mathbf{P}^0 and \mathbf{P}^* have the size, i.e. $d \times c$. That is, we can obtain only
 225 c projective vectors for feature extraction.

226 *C. The Review of LPP*

227 LPP [32], [33] computes the best linear approximations to
 228 the eigenfunctions of the manifold's Laplace Beltrami opera-
 229 tor. It aims to preserve local information and to find an embed-
 230 ding subspace which detects the most essential data manifold

231 structure [61], [62]. The objective function of LPP is to mini-
 232 mize

$$\begin{aligned} \frac{1}{2} \sum_{ij} (\mathbf{y}_i - \mathbf{y}_j)^2 \bar{w}_{ij} &= \frac{1}{2} \sum_{ij} \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|^2 \bar{w}_{ij} \\ &= \text{tr} (\mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B}) = \text{tr} (\mathbf{B}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{B}) \end{aligned} \quad (6)$$

233 where matrix \mathbf{B} denotes the transformation matrix, \mathbf{y}_i and \mathbf{y}_j
 234 denote the low-dimensional vectors of \mathbf{x}_i and \mathbf{x}_j in subspace
 235 \mathbf{B} , respectively. $\bar{\mathbf{W}}$ is supposed to be the similarity matrix of
 236 all pairwise data points, $\mathbf{L} = \bar{\mathbf{D}} - \bar{\mathbf{W}}$ is Laplacian matrix. $\bar{\mathbf{D}}$
 237 is a diagonal matrix and its element \bar{d}_{ii} is column or row sum
 238 of matrix $\bar{\mathbf{W}}$ (because $\bar{\mathbf{W}}$ is symmetric), i.e. $\bar{d}_{ii} = \sum_i \bar{w}_{ij}$.
 239 The similarity matrix \bar{w}_{ij} is defined as:

$$\bar{w}_{ij} = \begin{cases} \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/t), & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\|^2 < \varepsilon, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

240 where parameter $t \in R$, ε denotes the radius of the local neigh-
 241 borhood and it can be a sufficiently small positive value ($\varepsilon > 0$).
 242 In Eq. (7), the similarity matrix \bar{w}_{ij} might be sensitive to the
 243 value of the parameter t . To solve this problem, recently a
 244 parameter-free method was proposed in [63].

245 By considering the similarity matrix $\bar{\mathbf{W}}$, the relationship be-
 246 tween each data pair \mathbf{x}_i and \mathbf{x}_j in original space can be preserved
 247 by reconstructing the relationship between \mathbf{y}_i and \mathbf{y}_j in the low
 248 dimensional space \mathbf{B} with $\sum_{ij} (\mathbf{y}_i - \mathbf{y}_j)^2 \bar{w}_{ij}$ where $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i$
 249 and $\mathbf{y}_j = \mathbf{B}^T \mathbf{x}_j$. The optimal projections can be obtained by
 250 solving the following generalized eigen-function:

$$\mathbf{X} \mathbf{L} \mathbf{X}^T b = \lambda \mathbf{X} \bar{\mathbf{D}} \mathbf{X}^T. \quad (8)$$

251 Suppose $\lambda_i (i = 1, 2, \dots, d)$ are eigenvalues of problem 8,
 252 we can sort the eigenvalues in ascending order, then matrix
 253 $\mathbf{B} = [\mathbf{b}^1, \mathbf{b}^2, \dots, \mathbf{b}^k]$ combined of k eigenvectors corresponding
 254 to the first k smallest eigenvalues is the final projection matrix
 255 of LPP.

256 *III. THE EXTENSION BASED ON RIDGE REGRESSION*

257 In this section, we firstly review the definition of $L_{2,1}$ -norm
 258 and its property. Then we analyze the advantages and disadvan-
 259 tages of ridge regression. Meanwhile, we also propose a simple
 260 extension based on ridge regression.

261 *A. The Definition of $L_{2,1}$ -Norm and Its Property*

262 Some well-known models such as PCA, multilinear PCA
 263 (MPCA) [64], etc. use L_2 -norm as the measurement to com-
 264 pute the optimal projections in computer vision and face recog-
 265 nition. However, a large amount of experimental results have
 266 shown that in sparse feature selection, L_1 -norm outperforms
 267 L_2 -norm because of its generalization and the robustness for
 268 classification [16], [17], [61]. By combining the advantages of
 269 both L_1 -norm and part property of L_2 -norm, researchers obtain
 270 joint $L_{2,1}$ -norm minimization on both loss functions and regu-
 271 larization term for robust sparse learning for feature extraction
 272 [53]. Therefore, we use the $L_{2,1}$ -norm instead of L_2 -norm as a

273 new measurement for model design to overcome the problem of
 274 L_2 -norm being sensitive to outliers in a certain sense [54].

275 The $L_{2,1}$ -norm of a matrix is defined as

$$\|\mathbf{M}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^m m_{ij}^2} = \sum_{i=1}^n \|\mathbf{m}^i\|_2 \quad (9)$$

276 where the i -th row and the j -th column of a matrix $\mathbf{M} = (\mathbf{m}_i)_j$
 277 are denoted as \mathbf{m}^i and \mathbf{m}_j .

278 The common advantage of $L_{2,1}$ -norm and L_1 -norm based
 279 loss function is that they are more robust to outliers. However,
 280 the major difference between $L_{2,1}$ -norm and L_1 -norm is that
 281 $L_{2,1}$ -norm regularization is suitable for selecting meaningful or
 282 more powerful discriminant features from the data points with
 283 joint sparsity. The $L_{2,1}$ -norm based regularized methods can
 284 eliminate those useless interferences via making the elements
 285 in some rows of the projection matrix become zero such that
 286 the important features of the data points are emphasized and the
 287 insignificant features are ignored (filtered out) when conducting
 288 feature selection or extraction. Another advantage of $L_{2,1}$ -norm
 289 is that the $L_{2,1}$ -norm based methods are fast convergent and
 290 thus the computational cost is lower (this can be verified from
 291 computational cost of the $L_{2,1}$ -norm based methods compared
 292 with the L_1 -norm based methods in Table XI in Experiment
 293 section) [54], [56].

294 In all, employing the $L_{2,1}$ -norm instead of L_1 -norm as the
 295 regularization can obtain the joint sparsity to improve the per-
 296 formance and at the same time reduce the computational cost for
 297 efficient feature extraction and selection on image recognition
 298 tasks [54].

299 B. A Key Drawback in Traditional Regression

300 In (1) and (3), there exists a problem that when the number
 301 of the classes is too small, the traditional models cannot obtain
 302 enough projections for achieving good performance in pattern
 303 recognition. Thus, it is possible that learning more projection
 304 may improve the performance in feature extraction and classi-
 305 fication [4]. In order to obtain more projections in the regres-
 306 sion model, a tractable approach is to modify the representation
 307 $\|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2$ to be $\|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2$, which means that the
 308 matrix $(\mathbf{B} \mathbf{A}^T) \in R^{d \times c}$ takes the place of the matrix $\mathbf{P} \in R^{d \times c}$
 309 in the model. Thus we have the following optimization problem:

$$(\mathbf{A}^*, \mathbf{B}^*) = \arg \min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2, \text{ s.t. } \mathbf{A}^T \mathbf{A} = \mathbf{I}, \quad (10)$$

310 where \mathbf{A} is a $c \times k$ matrix and the size of matrix \mathbf{B} is $d \times k$
 311 where the notation k is any positive integer and c denotes the
 312 number of classes. In other words, the optimal solution \mathbf{B} with
 313 size $d \times k$ is able to break out the limitation of class number of
 314 the training data since the size of \mathbf{B} is not related to the class
 315 number and the variable k in \mathbf{B} is not related to the class number
 316 and the variable k in \mathbf{B} can be set as value that is larger than
 317 c , while \mathbf{P} with size $d \times c$ indicates that it can obtain at most c
 318 projections for feature selection.

From (10), we have

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{Y}) - \text{Tr}(2 \mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} - \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B}). \end{aligned} \quad (11)$$

319 By setting the derivatives of (11) with respect to \mathbf{B} equaling to
 320 0, the problem (11) is minimized at
 321

$$\mathbf{B}^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}, \quad (12)$$

322 where \mathbf{B}^* represents the optimal solution of (10).
 323

324 Denote the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, where $\mathbf{A}^T \mathbf{A} = \mathbf{I}$,
 325 $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, we have
 326

$$\mathbf{B}^* = \mathbf{U} \frac{1}{\mathbf{D}} \mathbf{V}^T \mathbf{Y} \mathbf{A}. \quad (13)$$

327 For (1) and 10), we have following propositions:
 328

329 *Proposition 1:* Suppose $\mathbf{X} \mathbf{X}^T$ is the full-rank matrix. Let
 330 \mathbf{P}^0 be the optimal solution to (1) and \mathbf{B}^* be the optimal solution
 331 to 10), if $k = c$ (i.e. the number of projection is equal to the
 332 number of class), then $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$.
 333

334 If $k > c$, the optimal solution \mathbf{B}^* with size $d \times k$ in (10) can
 335 obtain k projections instead of c projections as obtained by \mathbf{P}^0 in
 336 (1), which breaks out the small-class problem. In Proposition 1,
 337 the reason why the small-class problem is addressed by (10) is
 338 that the \mathbf{P}^0 in (1) has c projections while the optimal solution
 339 \mathbf{B}^* for (10) can learn k projections to perform feature extraction
 340 and classification, where k can be set as any integer. In other
 341 words, the number of the learned projections from (10) is not
 342 limited by the number of class and thus the small-class problem
 343 is addressed.

344 Similarly, for (3) and 10), we have following proposition:

345 *Proposition 2:* Let \mathbf{P}^* be the optimal solution to (3) and \mathbf{B}^*
 346 be the optimal solution to 10), if $k = c$, (i.e. the number of projec-
 347 tion is equal to the number of class), then $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$.
 348 Furthermore if $\alpha \rightarrow 0$, the metric matrices derived by \mathbf{B}^* and
 349 \mathbf{P}^* for classification are equivalent to each other.
 350

351 If $k > c$, the optimal solution \mathbf{B}^* with size $d \times k$ in (10) can
 352 obtain k projections instead of c projections as obtained by \mathbf{P}^*
 353 in (3), which breaks out the small-class problem.
 354

355 Proposition 2 indicates that when $\alpha \rightarrow 0$ (or using $\alpha = \varepsilon$,
 356 where ε is a very small number), the performance using \mathbf{B}^* and
 357 \mathbf{P}^* for feature extraction and classification will achieve the same
 358 results. If $k > c$, (10) can obtain more than c projections to per-
 359 form feature selection or extraction, which provides the theoreti-
 360 cal guarantee for the performance of (10). From Propositions 1
 361 and 2, we can draw the following conclusion:
 362

363 *Corollary 1:* If $\alpha \rightarrow 0$ and the matrix $\mathbf{X} \mathbf{X}^T$ is nonsingular,
 364 (1) and (3) have the same solution space.
 365

366 C. Other Drawbacks of Ridge Regression

367 Adding L_2 -norm term for regression is of great importance to
 368 deal with the singular problem in (1). Moreover, it shows that no
 369 matter the matrix $\mathbf{X} \mathbf{X}^T$ is singular or nonsingular, the classical
 370 ridge regression in (3) is able to obtain the optimal solution and
 371 (1) is only a special case of (3) with the regularization parameter
 372 $\alpha = 0$. However, there are still some obvious disadvantages in
 373

(3) since the optimal solution \mathbf{P}^* is not sparse, and thus it loses the feature selection function. Furthermore, the optimal solution \mathbf{P}^* in classical ridge regression model only contains the global information of the dataset and it ignores the local geometric structure. Thus, it is necessary to develop a new algorithm to deal with the above problems so as to enhance the effectiveness in feature extraction and pattern recognition. In the next section, we will propose a new model by jointing $L_{2,1}$ -norm regularization and locality regression to deal with the above problems.

IV. JOINTLY SPARSE LOCALITY REGRESSION ANALYSIS

In this section, the motivations and discussion are firstly present and then the proposed objective optimization problem as well as the local optimal solution will be presented.

A. The Motivations and Discussion

Based on the discussion in Section III-C and III-D, we can conclude the drawbacks of most existing regression methods into three aspects. First, due to the limitation of small-class problem, most regression methods cannot obtain enough projections to discover an effective projection matrix for discriminant feature extraction and classification. Second, the local structure of the data plays an important role in reconstructing the relationship between different data pairs in the low dimensional space. However, most existing regression methods do not take the local structure into consideration when performing feature selection or extraction. Third, there is no specific regression methods that are designed as regression form incorporating the local structures of the data as well as the sparsity of projections for feature selection and extraction.

Currently, deep learning technique is a research hotspot and it has been applied to the tasks of face recognition and object classification [65]. In spite of the high recognition rate of deep learning methods, behind is large-scale computing and long-term training. What is more, when the amount of data is not large enough, using deep learning methods for classification tends to obtain low performance because of the overfitting. In addition, most feature extraction methods based on deep learning [66], [67] do not consider the local structures of the data when doing convolutional operations. Even though they can obtain more abstract interpretation of the data, the relationship among different images is still missing. Therefore, developing efficient traditional feature extraction methods is still necessary for face recognition.

In conclusion, it is desirable to design a method that can solve the drawbacks of the existing regression methods and improve the performance of feature extraction to obtain high recognition rate with less computing time compared to the time-consuming and complicated deep learning methods.

B. The Objective Function of JSLR

To deal with the problems presented in Section III-C, Jointly Sparse Locality Regression Analysis (JSLR) is proposed to obtain a subset of jointly sparse projections for feature extraction and selection from the original data set. We also introduce the locality preserving regularized term to the model so as to characterize the local geometric structure of the data. Thus,

we present the objective function with joint $L_{2,1}$ -norm penalty and locality regularization. Let $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k]$ and $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k]$ be the variables of the following regression problem:

$$\begin{aligned} \bar{\mathbf{A}}, \bar{\mathbf{B}} = \arg \min_{\mathbf{A}, \mathbf{B}} & \left(\sum_{i=1}^n \|\mathbf{y}_i - \mathbf{x}_i^T \mathbf{B} \mathbf{A}^T\|_2^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \right. \\ & \left. + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|_2^2 \bar{\mathbf{w}}_{ij} \right) \\ & \mathbf{A}^T \mathbf{A} = \mathbf{I}, \end{aligned} \quad (14)$$

or in the matrix form

$$\begin{aligned} (\bar{\mathbf{A}}, \bar{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} & \left(\|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \right. \\ & \left. + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \bar{\mathbf{x}}_i - \mathbf{B}^T \bar{\mathbf{x}}_j\|_2^2 \bar{\mathbf{w}}_{ij} \right), \\ & \mathbf{A}^T \mathbf{A} = \mathbf{I}, \end{aligned} \quad (15)$$

where α_1 and α_2 are the regularization parameters. Since (14) and (15) have two variables and two kinds of norms in the model, they are not easy to be solved directly. Therefore, an alternatively iterative approach will be developed to solve the optimization problem in next section.

C. The Solutions of JSLR

From the definition of the $L_{2,1}$ -norm on the projection matrix \mathbf{B} , we have the diagonal matrix \mathbf{D}_B denoted as [54]

$$(\mathbf{D}_B)_{ii} = \frac{1}{2\|\mathbf{b}^i\|_2}, \quad (16)$$

where \mathbf{b}^i represents the i -th row of matrix \mathbf{B} .

Then from [56], we have the following equation:

$$\|\mathbf{B}\|_{2,1} = \text{Tr}(\mathbf{B}^T \mathbf{D}_B \mathbf{B}). \quad (17)$$

With the above preparation, we have

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \\ & + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|_2^2 \bar{\mathbf{w}}_{ij} \\ & = \text{Tr}(\mathbf{Y}^T \mathbf{Y} - 2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B} \\ & \quad + \alpha_1 (\mathbf{B}^T \mathbf{D}_B \mathbf{B}) + \alpha_2 \mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B}). \end{aligned} \quad (18)$$

Since the optimization problem has two variables, we need to fix one to compute the other. For fixed \mathbf{A} , by setting the derivatives of (18) with respect to \mathbf{B} equaling to 0, (18) is minimized by

$$\bar{\mathbf{B}} = (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} \quad (19)$$

Hence, when \mathbf{A} is fixed, the objective function of 14 or 15 is minimized at the local optimal solution \mathbf{B} . When fixing \mathbf{B} , $\text{Tr}(\mathbf{Y}^T \mathbf{Y} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B} + \alpha_1 (\mathbf{B}^T \mathbf{D}_B \mathbf{B}) + \alpha_2 \mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B})$ becomes a constant and thus it can be

441 ignored. In such case, the following maximization problem gives
 442 the optimal solution to (18):

$$\max_{\mathbf{A}} \text{Tr}(\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A}) \text{ s.t. } \mathbf{A}^T \mathbf{A} = \mathbf{I}. \quad (20)$$

443 Let $\bar{\mathbf{A}}$ be the optimization of (20). From the Theorem 4 in [4],
 444 we have

$$\bar{\mathbf{A}} = \mathbf{U} \mathbf{V}^T, \quad (21)$$

445 where \mathbf{U}, \mathbf{V} is the SVD decomposition value of $\mathbf{Y}^T \mathbf{X}^T \mathbf{B}$.

446 In addition, we can also have the following conclusion from
 447 above formulation:

448 *Theorem 1:* Let $\bar{\mathbf{B}}$ be the local optimal solution of the optimi-
 449 zation problem (14) or (15). If $\alpha_1 \rightarrow 0$ and $\alpha_2 \rightarrow 0$, the linear
 450 subspace spanned by the optimal solution of (14) or (15) approxi-
 451 mates to the linear subspace spanned by \mathbf{P}^0 and \mathbf{P}^* , namely,
 452 $\text{span}(\bar{\mathbf{B}}) = \text{span}(\mathbf{P}^0)$ and $\text{span}(\bar{\mathbf{B}}) = \text{span}(\mathbf{P}^*)$.

453 *Proof:* The proof is in the Appendix.

454 For (19), when $\alpha_1 = 0$ and $\alpha_2 = 0$, then $\bar{\mathbf{B}} = (\mathbf{X} \mathbf{X}^T)^{-1}$
 455 $\mathbf{X} \mathbf{Y} \mathbf{A} = \mathbf{B}^* = \mathbf{P}^0 \mathbf{A}$, where \mathbf{P}^0 and \mathbf{B}^* is the optimal solu-
 456 tion corresponding to (1) and (10). As Proposition 1 and
 457 Proposition 2 have presented the relationship between (10)
 458 and (1), (10) and (3) respectively, it is easy for us to have the
 459 following conclusions: ■

460 *Corollary 2:* With the same assumptions and notations as in
 461 Theorem 1, when $\alpha_2 = 0, \alpha_1 \rightarrow 0$, \mathbf{P}^0 and $\bar{\mathbf{B}}$ have the same
 462 linear subspace, namely, $\text{span}(\mathbf{P}^0) = \text{span}(\bar{\mathbf{B}})$.

463 *Corollary 3:* With the same assumptions and notations as in
 464 Theorem 1, when $\alpha_1 = 0, \alpha_2 \rightarrow 0$, \mathbf{P}^0 and $\bar{\mathbf{B}}$ have the same
 465 linear subspace, namely, $\text{span}(\mathbf{P}^0) = \text{span}(\bar{\mathbf{B}})$.

466 In summary, from Theorem 1, Corollary 2 and Corollary 3,
 467 we can know that either (14) or (15) provides a basic theoreti-
 468 cal guarantee for the effectiveness of the proposed regression
 469 model. Namely, when the parameters of the proposed model are
 470 set suitably, the optimal solution space of the ridge regression
 471 can be derived from (15). This means that the optimal projection
 472 of JSLR can approximate to the subspace spanned by the tradi-
 473 tional regression models. Besides, by utilizing the advantages
 474 of $L_{2,1}$ -norm regularization and locality preserving property, the
 475 proposed model is able to compute the jointly sparse projections
 476 and preserve the local geometric structure of the data for feature
 477 extraction. The detail of the iterative algorithm was illustrated
 478 in Algorithm 1.

479 D. Comparison and Discussion

480 In this section, we compare our algorithm JSLR with other
 481 methods, such as PCA, SPCA, LDA, LPP and so on. Both PCA
 482 and SPCA are outstanding in data processing and dimensionality
 483 reduction. PCA projects the original d -dimensional data onto
 484 $k (< d)$ -dimensional linear subspace with the combination of
 485 all the original variables. SPCA aims to produce modified sparse
 486 principal components by lasso (or elastic net) technique. But it
 487 just focuses on the global structure of the original data and ignore
 488 the local structure. Different from SPCA, JSLR can efficiently
 489 preserve the local geometric structure of the data set.

490 Some other subspace learning algorithms, LPP, NPE, etc. are
 491 able to preserve local structure of the original data. However,
 492 they cannot provide the jointly sparse property for the learned

Algorithm 1: JSLR Algorithm

Input: The training data $\mathbf{X} \in R^{d \times n}$,

the training data label $\mathbf{Y} \in R^{n \times c}$,

matrices $\bar{\mathbf{D}} \in R^{n \times n}$, $\bar{\mathbf{W}} \in R^{n \times n}$,

the objective dimension $k (k = 1, 2, \dots, n)$,

maximum number of the iteration: maxStep .

Step 1: Compute matrices $\bar{\mathbf{D}}, \bar{\mathbf{W}}$, and initialize matrix \mathbf{D}_B ,
 $step = 0$, $\text{converged} = \text{false}$.

Step 2: While !converged and $step <= \text{maxStep}$

- Compute \mathbf{B} using

$$\mathbf{B} = (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}$$

- Compute \mathbf{D}_B using $(\mathbf{D}_B)_{ii} = \frac{1}{2\|\mathbf{b}^i\|_2}$

- Compute \mathbf{A} using $\mathbf{A} = \mathbf{U} \mathbf{V}^T$

- Update $\text{converged} = \text{true}$ when \mathbf{B} is approximately
 changeless.

Step 3: Standardize the matrix \mathbf{B} to a final normalized
 matrix

and return it for feature selection.

Output: Low-dimensional discriminative subspace

$$\mathbf{B} \in R^{d \times k}, k = 1, 2, \dots, n.$$

493 subspace. Compared with them, JSLR achieves this goal by
 494 adding $L_{2,1}$ -norm regularization to make the elements in some
 495 rows of the projection to be 0 for efficient feature extraction and
 496 selection.

497 Ridge regression is frequently used in face recognition. How-
 498 ever, when the class number of training sample is too small, ridge
 499 regression cannot obtain more projections than the number of
 500 the classes for feature extraction. The same problem exists in
 501 LDA. In contrast, the number of the projections of JSLR is not
 502 limited by the number of the classes in training data. In spite
 503 of given a small number of classes in training sample set, JSLR
 504 can obtain any number of projections for feature selection and
 505 the number of the projections is freely set by the users.

506 In summary, the advantages of JSLR against PCA, SPCA,
 507 LDA, ridge regression and LPP are that JSLR can obtain joint
 508 sparsity and preserve the local structure for pattern recogni-
 509 tion. Another major difference between JSLR and other clas-
 510 sical methods is that the number of the training sample classes
 511 in JSLR is allowed to be very small but it still can learn more
 512 projections than the number of classes. These advantages make
 513 JSLR achieve high recognition rate.

V. THEORETICAL ANALYSIS

514 In this section, we present the theoretical analysis including
 515 convergence analysis and computational complexity analysis.

A. The Convergence

516 To verify the convergences of the proposed iterative algo-
 517 rithm, we begin with the following Lemmas:

518 *Lemma 1:* [54] For any two non-zero constants a and b , we
 519 have the following inequality:

$$\sqrt{a} - \frac{a}{2\sqrt{b}} \leq \sqrt{b} - \frac{b}{2\sqrt{b}}. \quad (22)$$

522 *Lemma 2*: [54] Denoted \mathbf{V} as any nonzero matrix, $\mathbf{V} \in R$,
 523 the following inequality holds:

$$\sum_i \|\mathbf{v}_t^i\|_2 - \sum_i \frac{\|\mathbf{v}_t^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2} \leq \sum_i \|\mathbf{v}_{t-1}^i\|_2 - \sum_i \frac{\|\mathbf{v}_{t-1}^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2}, \quad (23)$$

524 where $\mathbf{v}_t^i, \mathbf{v}_{t-1}^i$ denote the i -th row of matrix \mathbf{V}_t and \mathbf{V}_{t-1} .

525 *Proof*: Let $\|\mathbf{v}_t^i\|_2^2$ and $\|\mathbf{v}_{t-1}^i\|_2^2$ be the substitute of a and b
 526 in (22), the following inequality is valid for any i .

$$\|\mathbf{v}_t^i\|_2 - \frac{\|\mathbf{v}_t^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2} \leq \|\mathbf{v}_{t-1}^i\|_2 - \frac{\|\mathbf{v}_{t-1}^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2}, \quad (24)$$

527 Thus, (23) as the sum form of (24) also holds (22). With
 528 the above Lemma 1 and Lemma 2, we have the following
 529 theorem: ■

530 *Theorem 2*: Given all the parameters in the objective function
 531 except \mathbf{A} and \mathbf{B} , the iterative approach shown in Algorithm
 532 1 will monotonically decrease the objective function value of
 533 (14) or (15) in each iteration and provides a local optimal solution
 534 of the problem.

535 *Proof*: For simplicity, we denote the objective function of
 536 (18) as $F(\mathbf{B}, \mathbf{A}) = F(\mathbf{B}, \mathbf{A}, \mathbf{D}_B)$. Suppose for the $(t-1)$ -th
 537 iteration, both \mathbf{A}_{t-1} and \mathbf{B}_{t-1} can be obtained. Then we have
 538 the following inequality from (19):

$$F(\mathbf{B}_t, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}) \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (25)$$

539 For \mathbf{A}_t , as its optimal value comes from SVD and this will
 540 further decrease the value of the objective function, it goes

$$F(\mathbf{B}_t, \mathbf{A}_t, (\mathbf{D}_B)_{t-1}) \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (26)$$

541 In (18), since $\mathbf{Y}^T \mathbf{Y}$ is a constant, it can be ignored and we need
 542 to minimize

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}) \end{aligned}$$

543 As we have obtained the optimal \mathbf{B}_t and \mathbf{A}_t , then the following
 544 inequality holds:

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) \quad (27) \end{aligned}$$

545 That is

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) + \alpha_1 \sum_i \frac{\|\mathbf{b}_t^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) + \alpha_1 \sum_i \frac{\|\mathbf{b}_{t-1}^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \quad (28) \end{aligned}$$

TABLE I
THE COMPUTATIONAL COMPLEXITIES

Iteration variable	computational complexities
\mathbf{B}	$O(d^3)$
\mathbf{D}_B	$O(d^2)$
\mathbf{A}	$O(d^3)$

546 Then the above inequality indicates

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) \\ & + \alpha_1 \sum_i \|\mathbf{b}_t^i\|_2 - \alpha_1 \left(\sum_i \|\mathbf{b}_t^i\|_2 - \sum_i \frac{\|\mathbf{b}_t^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \right) \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) \\ & + \alpha_1 \sum_i \|\mathbf{b}_{t-1}^i\|_2 - \alpha_1 \left(\sum_i \|\mathbf{b}_{t-1}^i\|_2 - \sum_i \frac{\|\mathbf{b}_{t-1}^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \right) \quad (29) \end{aligned}$$

547 According to Lemma 2, we further have

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) + \alpha_1 \sum_i \|\mathbf{b}_t^i\|_2 \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) + \alpha_1 \sum_i \|\mathbf{b}_{t-1}^i\|_2. \quad (30) \end{aligned}$$

548 That is

$$\begin{aligned} & F(\mathbf{B}_t, \mathbf{A}_t) = F(\mathbf{B}_t, \mathbf{A}_t, (\mathbf{D}_B)_t) \\ & \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}) = F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (31) \end{aligned}$$

549 From (31), we can conclude that the objective function value
 550 of (14) or (15) is monotonically decreased via the updating rule
 551 presented in Algorithm 1. Therefore, the proposed iterative
 552 algorithm finally converges to the local optimal solution. ■

B. Computational Complexity Analysis

553 For simplicity, we assume the dimension of training samples
 554 is d . Our proposed algorithm aims to compute the matrix \mathbf{A} and
 555 \mathbf{B} . Computing \mathbf{B} in (19) needs $O(d^3)$ while computing \mathbf{D}_B
 556 in (16) needs $O(d^2)$. Since SVD of $\mathbf{Y}^T \mathbf{X}^T \mathbf{B}$ also needs $O(d^3)$,
 557 then the computational complexity of \mathbf{A} is also $O(d^3)$. It is easy
 558 to know that the main complexity of the algorithm is $O(Td^3)$,
 559 where T denotes the number of iterations for convergence. Ta-
 560 ble I lists the computational complexities of each variable.

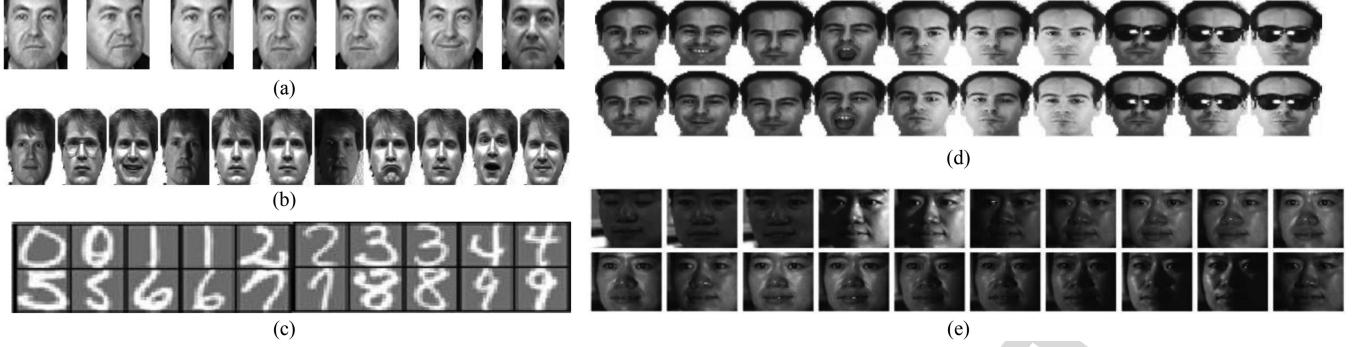


Fig. 1. Examples from FERET, Yale, USPS, AR, and CMU PIE data sets. (a) FERET. (b) Yale. (c) USPS. (d) AR. (e) CMU PIE.

VI. EXPERIMENTS

In this section, to evaluate the proposed JSLR algorithm for feature selection, we conducted a set of experiments from three aspects: experiments on small-scale face databases, experiments on large-scale databases and experiments based on deep learning. In experiments, several classical as well as state-of-the-art methods are used as compared methods. They are the classical principle component analysis method PCA, the classical Ridge Regression (RR) [60], the Linear Discriminant Analysis (LDA) [5], the traditional sparse learning method SLDA based on L_1 -norm [23], the local structure learning method Locality Preserving Projection (LPP) [32], the regression analysis of locality preserving projections via sparse penalty (SpLPP) [61] which applies sparsity penalty and minimization based on L_1 -norm to locality preserving projections, the dictionary learning methods (i.e. label consistent K-SVD (LC-KSVD2) [68] and the Locality Constrained and Label Embedding Dictionary Learning (LCLE-DL) [69]), the most related $L_{2,1}$ -norm regularization methods for feature selection and subspace learning (i.e. Unsupervised Discriminative Feature Selection (UDFS) [56] and Robust Feature Selection (RFS) [54]). In addition, the proposed method without $L_{2,1}$ -norm regularization named JSLR($\alpha_1 = 0$) (i.e. the second term in the proposed objective function in Eq. (14) is removed) was added as a compared method to all experiments to evaluate the effectiveness of the jointly sparse regularization.

In all experiments we make comparison in the avenue of deep learning (the method is called Deep-NN in this paper). Deep-NN is completed by the following two steps. Firstly, we use the deep convolutional neural network (CNN) as the feature extractor to obtain the deep features of all samples. This process is similar to [70]. Secondly, we use the nearest neighbor classifier (NN) for classification. For the proposed JSLR, we also use the deep features instead of the traditional image features as input and this method is called Deep-JSLR for easy understanding. Note that the deep features of character database are obtained according to the tutorial of MNIST network on official Caffe site (<http://caffe.berkeleyvision.org/gathered/examples/mnist.html>.)

A. Experiments on Small-Scale Database

In this section, experiments on four databases, including FERET, AR, CMU PIE and Yale database, were conducted to

evaluate the performance of the proposed method versus the compared methods under different variations of facial expression and lighting condition.

1) *Experiments on FERET Face Database*: The FERET face database [71] includes 1,400 images of 200 individuals (each individual has seven images). In the experiment, the facial portion of each original image was automatically cropped based on the location of the eyes, and the cropped images were resized to 40×40 pixels. The sample images of one person are shown in Fig. 1(a).

Experimental Setting: For all the databases, the image set is partitioned into two parts, i.e. the gallery and probe sets. In each database, l (l is no more than the number of class) images of each class are randomly selected to form the gallery set and the remaining images are used as the probe set. PCA was used as pre-processing to reduce the dimension of data. Then the proposed method and the compared methods were used to perform feature extraction, independently. Finally, nearest neighbor classifier was used for classification. The experiments were independently performed 10 times. The average recognition rates and the corresponding dimensions as well as the standard deviations of each method were listed on the Table II. Besides, the comparison results were also shown in the Fig. 2(c)–(f) when 5 images of each individual were randomly selected for training and the remaining images were used for testing. The dimensions of the projection matrices were set as empirical value and marked on the horizontal axis. The variables except parameter α_1 and α_2 in JSLR were randomly initialized in our experiments.

Exploration of the Performance of the Parameters: In order to explore the optimal parameters for JSLR on different data sets, we analyzed the values of the parameters α_1 (Alpha1) and α_2 (Alpha2). For the other compared methods, since in most of cases the best performance lie on the area of $[10^{-3}, 10^3]$, as introduced in the corresponding papers, we fixed their parameters on the area of $[10^{-3}, 10^3]$ and report the best results.

In this experiment, we analyze the impacts of various parameter values on the performance of JSLR and the average recognition rates of different dimensions from 5 to 200. Table II shows the best average recognition rates based on 10 times running and the corresponding dimensions as well as the standard deviations of each method with l ($l = 4, 5$) images of each individual for training while the remaining images were used for testing.

TABLE II
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF DIFFERENT METHODS ON FERET FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	54.55	54.72	60.77	46.80	61.23	53.12	63.82	55.53	59.42	46.72	73.22	74.15	99.47	100.00
	± 8.54	± 8.74	± 20.04	± 9.61	± 19.96	± 12.55	± 21.07	± 23.33	± 13.46	± 9.55	± 18.77	± 18.34	± 0.46	± 0.00
	$28*5$	$27*5$	$30*5$	$30*5$	$30*5$	$12*5$	$40*5$	$13*5$	$30*5$	$30*5$	$30*5$	$29*5$	$30*5$	$26*5$
5	65.50	65.90	76.45	61.50	77.18	66.65	78.35	74.92	67.95	60.08	90.70	91.35	99.55	100.00
	± 5.48	± 5.65	± 7.76	± 4.21	± 7.61	± 5.27	± 7.33	± 10.55	± 6.33	± 6.69	± 6.67	± 6.42	± 0.33	± 0.00
	$30*5$	$27*5$	$30*5$	$30*5$	$29*5$	$11*5$	$40*5$	$8*5$	$29*5$	$30*5$	$27*5$	$29*5$	$30*5$	$20*5$

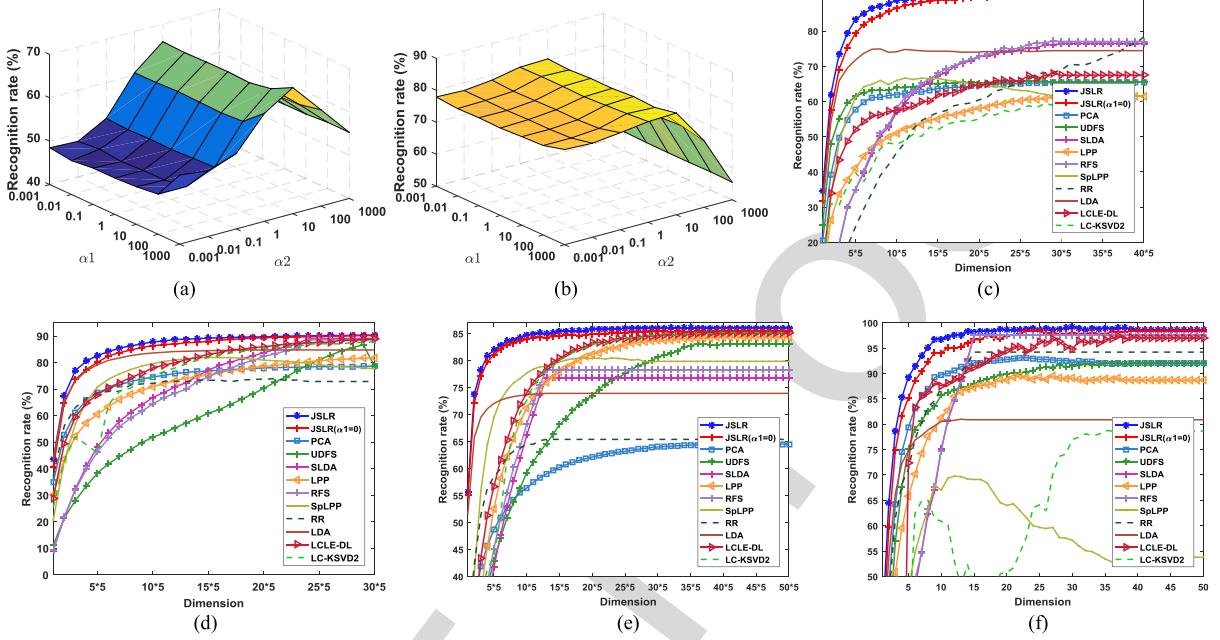


Fig. 2. The recognition rate versus the parameters α_1 and α_2 on the (a) FERET and (b) AR face database, respectively. The recognition rates (%) versus the dimensions of different methods on the (c) FERET, (d) AR, (e) CMU PIE, (f) Yale face databases, respectively.

Fig. 2(a) shows the recognition rates when the two parameters α_1 and α_2 change from 10^{-3} to 10^3 . Fig. 2(c) shows the average recognition rates versus various dimensions of different methods.

It is easy to know that the optimal value of the parameter α_1 lies on the area of $[10^{-3}, 10^2]$ while the optimal value of the parameter α_2 lies on the area of $[10^{-2}, 10^3]$. In other words, JSLR is efficient and robust among these areas. By contrast, when the values of the two parameters lie on other areas, it will cause the larger decline of the recognition rates.

As it can be seen from Fig. 2(c), the recognition rates of JSLR as well as $\text{JSLR}(\alpha_1 = 0)$ are the highest. The results shown in Table II and Fig. 2(c) indicate that JSLR and $\text{JSLR}(\alpha_1 = 0)$ outperform PCA, SLDA, UDFS, RFS, LPP and SpLPP, RR, LDA, LCLE-DL, LC-KSVD2 in feature extraction. Besides, from Table II, we can easily know that Deep-JSLR outperforms Deep-NN.

2) *Experiments on AR Face Database*: The AR face database [72] contains the pictures of 120 individuals (each individual has 20 images). The face portion of each image was manually cropped (because of missing eye coordinates) and then normalized to 50×40 pixels. The sample images of one person are shown in Fig. 1(d).

In this experiment, we randomly selected l ($l = 4, 5, 6$) images of each individual for training, and the rest of the images in the data set were used for testing. From Fig. 2(b), we can know that the optimal values of parameter α_1 and α_2 were both $[10^{-3}, 10^2]$. Thus, we used this area for JSLR to obtain the comparison results. Table III listed the performance of different methods. Fig. 2(d) showed the average testing recognition rates. It is obvious that JSLR or Deep-JSLR outperforms the other methods.

3) *Experiments on CMU PIE Database*: The CMU PIE face database [73] contains 68 individuals with 41,368 face images as a whole. We selected a subset (C29) containing 1632 images from 68 individuals (each providing 24 images). All of these face images were automatically aligned based one-eye coordinates and cropped to 32×32 pixels. Fig. 1(e) shows the sample images from this database.

In this experiment, l ($l = 4, 5, 6$) images of each individual were randomly selected for training, and the rest of the images in the data set were used for testing. The optimal areas of α_1 and α_2 were the same with the areas on AR database. Table IV presents the performance of different methods. Fig. 2(e) shows the average testing recognition rates and indicates that JSLR outperforms the other methods again.

TABLE III
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON AR FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	76.86	85.66	87.97	80.40	87.41	79.03	68.29	80.06	89.89	84.97	88.82	89.08	83.43	90.23
	± 4.56	± 8.48	± 10.95	± 9.25	± 11.41	± 7.76	± 5.89	± 11.78	± 7.76	± 6.96	± 10.79	± 10.64	± 12.78	± 9.77
	30*5	29*5	24*5	30*5	24*5	12*5	21*5	23*5	29*5	30*5	30*5	28*5	29*5	28*5
5	78.67	86.73	89.07	81.91	88.58	82.01	73.91	84.78	88.87	87.88	90.19	90.33	89.26	94.70
	± 5.41	± 8.89	± 11.27	± 9.78	± 11.67	± 8.34	± 7.99	± 11.77	± 8.28	± 7.44	± 10.63	± 10.60	± 9.36	± 7.02
	30*5	29*5	24*5	30*5	24*5	15*5	21*5	21*5	30*5	30*5	29*5	29*5	29*5	29*5
6	80.47	90.35	94.24	86.11	93.85	86.76	79.41	91.82	91.50	92.26	95.05	95.23	93.27	97.52
	± 4.97	± 5.85	± 8.11	± 7.89	± 8.51	± 6.15	± 6.02	± 9.38	± 5.92	± 5.85	± 8.02	± 8.04	± 3.07	± 1.84
	30*5	28*5	24*5	30*5	24*5	16*5	21*5	22*5	30*5	30*5	28*5	29*5	29*5	30*5
15	92.08	92.08	99.22	98.80	99.22	95.03	96.30	99.12	98.02	97.10	99.55	99.60	97.68	99.73
	± 18.01	± 3.44	± 20.38	± 0.71	± 0.41	± 2.54	± 1.57	± 0.68	± 1.16	± 1.26	± 0.32	± 0.32	± 2.91	± 0.74
	30*5	28*5	24*5	30*5	24*5	16*5	21*5	21*5	30*5	30*5	19*5	30*5	30*5	30*5
16	92.35	92.35	99.27	98.92	99.27	95.19	96.69	99.25	98.08	97.50	99.63	99.77	98.67	99.98
	± 18.07	± 3.48	± 20.41	± 0.78	± 0.45	± 1.72	± 1.90	± 0.53	± 1.17	± 1.01	± 0.28	± 0.25	± 1.38	± 0.07
	30*5	28*5	24*5	30*5	24*5	16*5	21*5	21*5	30*5	30*5	18*5	19*5	30*5	22*5

TABLE IV
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON CMU PIE FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	57.11	77.62	69.25	80.01	71.55	76.21	56.33	36.31	81.43	81.05	80.60	81.43	96.32	97.54
	± 12.66	± 8.85	± 13.46	± 9.16	± 12.92	± 8.32	± 10.50	± 8.54	± 6.78	± 8.15	± 11.12	± 10.4	± 1.10	± 0.92
	40*5	37*5	13*5	40*5	13*5	20*5	13*5	8*5	39*5	36*5	34*5	30*5	38*5	37*5
5	64.52	83.13	76.76	84.35	78.25	80.48	65.39	73.92	85.11	85.22	85.44	86.08	96.73	97.96
	± 12.27	± 4.65	± 11.02	± 4.85	± 10.19	± 6.86	± 9.26	± 10.23	± 4.32	± 3.64	± 5.19	± 4.54	± 1.25	± 1.17
	40*5	36*5	13*5	40*5	13*5	19*5	13*5	13*5	37*5	39*5	33*5	34*5	39*5	34*5
6	68.84	85.63	78.45	85.62	79.77	83.51	71.36	78.19	85.50	85.28	85.16	85.85	96.99	98.14
	± 10.20	± 4.74	± 9.59	± 5.04	± 7.85	± 5.44	± 10.07	± 10.30	± 3.56	± 3.80	± 5.15	± 4.54	± 1.10	± 1.07
	40*5	38*5	13*5	37*5	13*5	19*5	13*5	13*5	40*5	35*5	38*5	33*5	39*5	38*5
19	94.24	94.24	92.29	92.09	92.09	91.88	94.41	91.06	87.00	90.09	92.91	93.12	99.32	100.00
	± 5.91	± 5.91	± 8.22	± 8.03	± 8.34	± 8.39	± 5.46	± 9.55	± 9.74	± 7.52	± 7.65	± 7.55	± 0.77	± 0.00
	40*5	38*5	13*5	40*5	13*5	19*5	13*5	13*5	40*5	40*5	40*5	33*5	14*5	
20	93.16	93.16	90.29	90.18	90.11	90.22	93.68	89.04	85.57	88.54	91.29	91.47	99.34	100.00
	± 7.17	± 7.17	± 10.25	± 9.95	± 10.28	± 9.95	± 6.56	± 11.73	± 11.96	± 9.17	± 9.44	± 9.28	± 0.76	± 0.00
	40*5	38*5	13*5	40*5	13*5	19*5	13*5	13*5	40*5	40*5	23*5	30*5	33*5	14*5

TABLE V
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON YALE FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	91.24	90.48	96.95	85.05	96.48	68.10	91.62	66.76	97.90	78.86	97.71	98.95	99.71	100.00
	± 3.46	± 4.02	± 2.49	± 5.06	± 2.58	± 9.27	± 3.24	± 10.23	± 1.41	± 41.68	± 1.61	± 1.05	± 0.40	± 0.00
	21	39	15	37	15	11	15	10	40	33	31	27	40	12
5	93.11	92.11	97.56	89.56	97.56	69.78	94.22	81.00	97.44	78.78	98.44	99.22	99.56	100.00
	± 3.60	± 4.57	± 2.39	± 3.58	± 2.21	± 5.49	± 3.57	± 5.50	± 1.80	± 41.78	± 1.57	± 1.20	± 0.47	± 0.00
	22	36	15	22	15	12	15	13	37	35	39	30	40	11

692 4) *Experiments on Yale Database*: The Yale face database
693 [43] contains 165 grayscale images of 15 individuals. Each
694 image was manually cropped (because of no eye coordinates
695 provided) and resized to 50×40 pixels. Fig. 1(b) shows the
696 sample images from this database.

697 In this experiment, l ($l = 4, 5$) images of each individual were
698 randomly selected for training, and the rest of the images in the
699 data set were used for testing. The values of α_1 and α_2 were both
700 from 10^{-3} to 10^2 . The performances of the different methods
701 are shown in Table V. Fig. 2(f) shows the average recognition
702 rates. It clearly indicates that JSLR and Deep-JSLR can obtain
703 the best performance when the traditional image features and
704 deep features are used as input.

705 B. Experiments on Large-Scale Database

706 In this section, two different databases are used to evaluate
707 the performance of the proposed method based on large-scale
708 data learning.

709 1) *Experiments on USPS Database*: The United States
710 Postal Service(USPS) database [55] consists of 1,100 of each
711 handwritten digit (0-9). The images in this database are resized
712 to 16×16 pixels. The performance of JSLR on large-sample
713 data set was evaluated on this database. Fig. 1(c) shows the sample
714 images from this database.

715 In this experiment, l ($l = 400, 500, 600, 5, 10$) images of each
716 class were randomly selected for training, and the rest of the images
717 in the data set were used for testing. The parameter α_1 and
718 α_2 were both from 10^{-3} to 10^3 . Table VI shows the performance
719 of the different methods. From the result, we can know that JSLR
720 or Deep-JSLR can achieve better performance than other compared
721 methods on this database. Particularly, when only 5/1100 images
722 of each class are used for training, the proposed method
723 can obtain higher recognition rate than not only the compared
724 methods but also the Deep-NN.

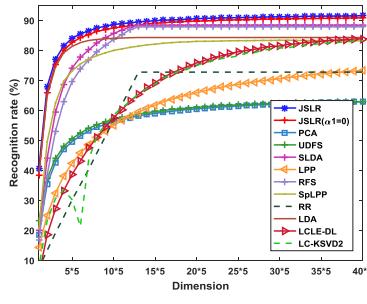
725 2) *Experiments on PIE67 \times 170 Database*: The PIE67 \times
726 170 database is a subset of the CMU PIE face database [73].

TABLE VI
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON USPS FACE DATABASE

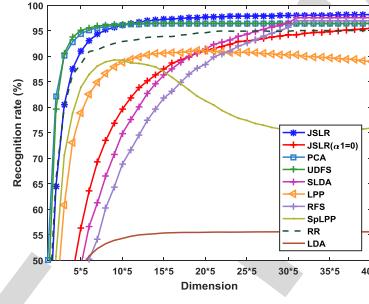
Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR	Deep-NN	Deep-JSLR
400	90.55 ± 0.13	91.17 ± 0.12	91.37 ± 0.35	89.41 ± 0.30	90.42 ± 0.40	91.02 ± 0.19	90.34 ± 0.17	86.78 ± 0.00	87.21 ± 0.56	90.78 ± 0.19	92.02 ± 0.22	92.35 ±0.18	91.44 ± 0.28	95.74 ±0.27
	7*5 $10*5$	2*5 $2*5$	6*5 $6*5$	6*5 $2*5$	6*5 $2*5$	6*5 $2*5$	6*5 $2*5$	1*5 $1*5$	4*5 $4*5$	30*5 $30*5$	30*5 $30*5$	20*5 1*5	1*5 $1*5$	15*5 21*5
	91.02 ± 0.15	91.75 ± 0.18	91.91 ± 0.23	90.14 ± 0.21	90.81 ± 0.31	91.90 ± 0.16	90.87 ± 0.23	87.23 ± 0.00	87.22 ± 0.59	91.60 ± 0.27	92.52 ± 0.18	92.77 ±0.14	92.18 ± 0.17	96.00 ±0.26
500	91.49 ± 0.17	92.20 ± 0.24	92.28 ± 0.35	90.58 ± 0.17	91.40 ± 0.45	92.51 ± 0.13	91.45 ± 0.25	87.49 ± 0.00	87.15 ± 0.71	92.15 ± 0.33	92.82 ± 0.28	93.08 ±0.31	92.67 ± 0.26	96.09 ±0.21
	7*5 $8*5$	2*5 $2*5$	2*5 $6*5$	6*5 $2*5$	6*5 $6*5$	6*5 $2*5$	6*5 $2*5$	1*5 $1*5$	4*5 $4*5$	30*5 $30*5$	28*5 $28*5$	25*5 1*5	1*5 $1*5$	20*5 21*5
	61.43 ± 2.08	61.43 ± 2.08	56.83 ± 2.30	49.41 ± 3.10	56.23 ± 2.43	15.76 ± 2.28	57.19 ± 2.24	61.45 ± 3.53	62.00 ± 2.40	62.04 ± 2.34	59.51 ± 2.24	67.83 ±2.39	56.96 ± 2.28	66.96 ±1.72
600	69.86 ± 2.18	69.86 ± 2.18	55.78 ± 2.15	49.49 ± 2.64	54.73 ± 2.00	14.14 ± 2.30	55.64 ± 2.05	61.02 ± 2.16	68.87 ± 1.75	67.81 ± 1.91	58.67 ± 1.87	76.22 ±2.64	65.66 ± 1.68	78.15 ±1.64
	7*5 $8*5$	8*5 $2*5$	2*5 $6*5$	6*5 $2*5$	6*5 $6*5$	2*5 $2*5$	1*5 $1*5$	4*5 $4*5$	30*5 $30*5$	12*5 $12*5$	12*5 $100*5$	20*5 100*5	100*5 $100*5$	20*5 20*5
	69.86 ± 2.18	69.86 ± 2.18	55.78 ± 2.15	49.49 ± 2.64	54.73 ± 2.00	14.14 ± 2.30	55.64 ± 2.05	61.02 ± 2.16	68.87 ± 1.75	67.81 ± 1.91	58.67 ± 1.87	76.22 ±2.64	65.66 ± 1.68	78.15 ±1.64



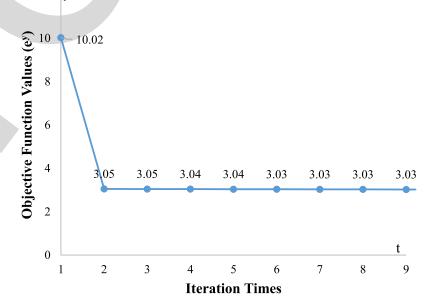
(a)



(b)



(c)



(d)

Fig. 3. (a) Sample images on the LFW database. The recognition rates (%) versus the dimensions of different methods on the (b) PIE67 \times 170, (c) LFW databases, respectively. (d) An example of the convergence curve of JSLR on Yale database.

TABLE VII
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON PIE67 \times 170 FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR	Deep-NN	Deep-JSLR
10	45.75 ± 3.71	45.75 ± 3.72	78.63 ± 3.55	59.72 ± 3.33	76.14 ± 4.03	69.71 ± 3.69	58.85 ± 3.01	72.74 ± 2.70	82.72 ± 2.86	78.78 ± 2.11	81.47 ± 2.93	83.64 ±3.17	90.90 ± 1.21	95.62 ±1.31
	40*5 $40*5$	40*5 $40*5$	13*5 $13*5$	40*5 $40*5$	13*5 $13*5$	21*5 $21*5$	13*5 $13*5$	13*5 $13*5$	40*5 $40*5$	40*5 $40*5$	40*5 $40*5$	40*5 40*5	40*5 $40*5$	40*5 40*5
	62.89 ± 3.05	62.89 ± 3.06	88.41 ± 2.09	73.32 ± 3.12	87.84 ± 2.11	83.38 ± 3.09	72.75 ± 2.66	84.29 ± 2.11	86.84 ± 1.42	86.19 ± 1.23	90.84 ± 1.42	91.55 ±1.37	93.98 ± 1.07	97.70 ±0.73
20	69.85 ± 3.57	69.85 ± 3.58	91.46 ± 2.22	77.99 ± 2.76	91.05 ± 2.30	87.42 ± 2.72	83.74 ± 3.13	87.98 ± 2.78	87.81 ± 1.61	87.96 ± 1.27	93.38 ± 1.82	93.70 ±1.78	95.36 ± 1.00	98.38 ±0.67
	40*5 $40*5$	40*5 $40*5$	13*5 $13*5$	40*5 $40*5$	13*5 $13*5$	36*5 $36*5$	13*5 $13*5$	13*5 $13*5$	40*5 $40*5$	39*5 $39*5$	40*5 $40*5$	40*5 40*5	40*5 $40*5$	40*5 40*5
	69.85 ± 3.57	69.85 ± 3.58	91.46 ± 2.22	77.99 ± 2.76	91.05 ± 2.30	87.42 ± 2.72	83.74 ± 3.13	87.98 ± 2.78	87.81 ± 1.61	87.96 ± 1.27	93.38 ± 1.82	93.70 ±1.78	95.36 ± 1.00	98.38 ±0.67

There are total 11,390 images from 67 individuals and each individual has 170 images on this database. The experiment on this database is conducted to evaluate the performance of JSLR as well as the compared methods on the occasion when there are various facial expression, lighting condition and angle on the face images.

In this experiments, l ($l = 10, 20, 30$) images of each individual are randomly selected for training and the remaining are used for testing. The recognition rates of all methods are shown

in Fig. 3(b) and Table VII. From Fig. 3(b), JSLR as well as JSLR($\alpha=0$) obtain higher recognition rate than other methods, which indicates that the proposed method is superior to other methods even without the $L_{2,1}$ -norm regularization term (this can also be verified by Fig. 2(c) and Table II).

C. Experiments Based on Deep Learning

In this section, experiments on three database (AR, the standard subsets of the FERET and the LFW databases [74]) were

TABLE VIII
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON AR FACE DATABASE BASED ON DEEP LEARNING

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	Deep-NN	JSLR ($\alpha=0$)	JSLR
4	87.12	87.12	87.53	83.46	88.80	87.12	74.53	82.37	90.01	89.89	83.43	13.49	90.23
	± 12.26	± 12.33	± 9.26	± 10.96	± 10.73	± 11.90	± 9.91	± 8.81	± 6.50	± 6.89	± 12.78	± 7.47	± 9.77
	30*5	30*5	24*5	22*5	24*5	30*5	24*5	23*5	30*5	30*5	29*5	30*5	28*5
5	91.84	91.84	91.16	89.82	93.17	91.84	80.74	89.63	94.67	94.10	89.26	76.91	94.70
	± 9.36	± 9.43	± 7.84	± 8.20	± 7.79	± 8.83	± 7.82	± 7.56	± 4.61	± 4.60	± 9.36	± 5.51	± 7.02
	30*5	30*5	24*5	20*5	24*5	30*5	24*5	21*5	30*5	30*5	29*5	30*5	29*5
6	95.40	95.33	94.39	93.99	96.16	95.33	84.64	93.85	96.81	96.67	93.27	92.93	97.52
	± 2.77	± 2.75	± 3.84	± 3.14	± 2.74	± 2.82	± 6.58	± 3.65	± 2.25	± 1.78	± 3.07	± 3.58	± 1.84
	29*5	30*5	24*5	16*5	24*5	30*5	24*5	18*5	30*5	30*5	29*5	30*5	30*5

TABLE IX
THE MAXIMAL RECOGNITION RATE OF ALL METHODS ON THE Fb, Fc, Dup1, Dup2 FACE DATABASE BASED ON DEEP LEARNING

Algorithm	Fb (dim=216, 512)	Fc (dim=216, 512)	Dup1(dim=216, 512)	Dup2(dim=216, 512)
PCA	99.41, 99.41	99.48, 99.48	98.20, 98.20	98.29, 98.29
UDFS	80.75, 99.41	65.46, 99.48	43.07, 98.20	45.30, 98.29
SLDA	98.49, 98.91	100.00, 100.00	84.63, 88.50	89.74, 92.31
LPP	98.58, 99.16	90.03, 92.66	86.75, 91.45	86.75, 91.45
RFS	99.58, 99.58	100.00, 100.00	97.51, 98.48	98.72, 99.15
SpLPP	99.41, 99.50	98.97, 99.48	97.78, 98.06	98.72, 99.15
RR	98.74, 98.91	98.45, 98.45	94.46, 94.46	96.58, 96.58
LDA	99.41, 99.41	99.48, 99.48	98.20, 98.20	98.29, 98.29
LCLE-DL*	-	-	-	-
LC-KSVD2*	-	-	-	-
Deep-NN	99.41, 99.50	99.48, 99.48	98.48, 98.48	98.72, 98.72
JSLR ($\alpha=0$)	98.58, 99.41	99.48, 100.00	86.01, 93.77	88.03, 94.44
JSLR	99.67, 99.67	100.00, 100.00	98.75, 98.89	99.57, 99.57

*Since there is only one sample in each class on the training set of Fa, the dictionary learning methods are not suitable to use in this case and the performance is too poor to be presented.

conducted based on deep learning. In the experiments, the Caffe deep learning framework [75] was used as the pre-processing to learn the deep features from the sample images. After the deep features were obtained, we further used the subspace learning methods (i.e. PCA, UDFS, SLDA, LPP, RFS, SpLpp, RR, LDA and the proposed JSLR) and dictionary learning methods (i.e. LCLE-DL and LC-KSVD2) to perform further feature extraction and then the nearest neighbor classifier was used for classification.

For AR and the standard subsets of the FERET databases, the dimension of extracted features based on the deep convolutional neural network (CNN) is 512 while that of the LFW database is 1024. For the standard FERET dataset, the Fa subset was used as the gallery set while the Fb, Fc, Dup1 and Dup2 were used as the probe sets. The LFW database contains images from 5,749 subjects in the uncontrolled environment, which makes it as a challenging recognition task. 158 subjects with total 4,324 images are selected from LFW-a subset and used in our experiment as the LFW-a subset is the aligned version of LFW database. The sample images on this database are shown in Fig. 3(a).

The experimental results on AR databases is listed in Table VIII. For the standard subsets of the FERET database, the best recognition rates corresponding different methods are

shown in Table IX, in which the accuracy on both the original dimensions (i.e. 512) and 216 dimensions (i.e. half of 512) are listed. The results in Table VIII and IX clearly show that the performance of the proposed JSLR is better than that of Deep-NN. This indicates that JSLR is able to extract discriminative information from deep features and further achieve higher recognition rate. The experimental results on LFW database are shown in Fig. 3(c) and Table X. In Fig. 3(c), the reason why no curves of LCLE-DL and LC-KSVD2 present is that no PCA is used as pre-processing to reduce the dimension of the input data to a specific value, the dictionary learning methods (i.e. LCLE-DL and LC-KSVD2) can only obtain recognition rate corresponding to the original dimension (i.e. 1024). Therefore, we cannot obtain the recognition rate curve versus the dimension variations for the two methods. From Fig. 3(c) and Table X, we can know that JSLR outperforms other compared methods again.

The convergence curves of the proposed JSLR on all databases are shown in Fig. 3(d) and Fig. 4. In these figures, the objective function value corresponding to each iteration is denoted as e^y where y is the values marked on the vertical coordinate. The convergence curves on all databases indicate that the proposed method can converge after several iterations.

TABLE X
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON THE LFW FACE DATABASE BASED ON DEEP LEARNING

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	Deep-NN	JSLR ($\alpha_1=0$)	JSLR
3	96.46	96.50	97.59	91.13	96.96	89.29	95.06	55.60	95.39	93.94	96.45	95.51	98.19
	± 0.18	± 0.18	± 0.23	± 0.67	± 0.19	± 1.02	± 0.43	± 3.36	± 0.57	± 0.96	± 0.16	± 0.39	± 0.20
	15*5	16*5	31*5	19*5	31*5	9*5	31*5	29*5	40*5	40*5	197*5	40*5	40*5
5	96.79	96.80	98.58	97.82	98.39	96.67	96.48	69.64	97.51	95.48	96.78	98.39	98.71
	± 0.23	± 0.22	± 0.22	± 0.20	± 0.26	± 0.24	± 0.32	± 2.22	± 0.32	± 0.30	± 0.20	± 0.14	± 0.20
	14*5	13*5	31*5	21*5	31*5	9*5	31*5	25*5	40*5	40*5	198*5	40*5	39*5
7	97.06	97.12	99.01	98.85	98.96	98.10	97.07	70.47	97.95	97.16	97.08	98.91	98.92
	± 0.11	± 0.15	± 0.14	± 0.17	± 0.10	± 0.20	± 0.32	± 2.71	± 0.35	± 0.38	± 0.12	± 0.14	± 0.13
	20*5	13*5	31*5	25*5	31*5	9*5	25*5	26*5	40*5	40*5	197*5	39*5	39*5

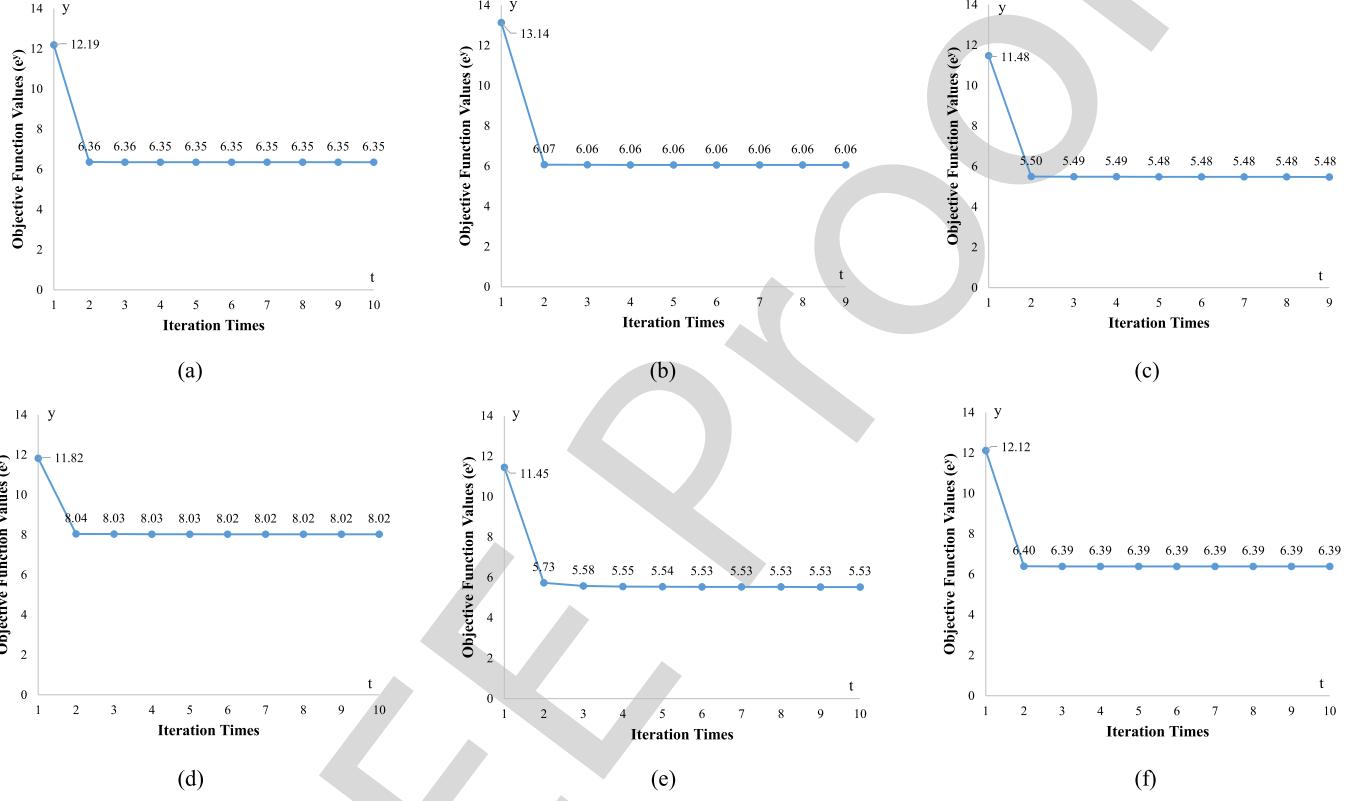


Fig. 4. Examples of the convergence curves of JSLR on (a) FERET, (b) AR, (c) CMU PIE, (d) USPS, (e) PIE67 \times 170 and (f) LFW database, respectively.

790 D. Experimental Results and Discussions

791 The comparison among the proposed JSLR, classical PCA,
792 RR, LDA, SLDA, LPP, $L_{2,1}$ -norm based methods
793 (UDFS, RFS) and dictionary learning methods (LCLE-DL, LC-
794 KSVD2) has been presented using recognition rates on these
795 databases: FERET, AR, CMU PIE, Yale and USPS. From the
796 results, we reveal the following interesting points:
797

- 798 1) In all experiments, including the face databases (FERET,
799 AR, CMU PIE, Yale) and non-facial database (USPS),
800 JSLR consistently achieves higher recognition rates than
801 other methods. These results are in line with the theoretical
802 analysis of JSLR that it obtains discriminative information
803 with joint sparsity and takes local geometric structure of
804 dataset into consideration to perform feature selection and
805 extraction.
806 2) JSLR is able to encode more discriminating information
807 in the low-dimensional face subspace since the local

808 geometric structure is considered to be more effective than
809 the global structure for feature extraction and feature se-
810 lection in some cases. The reason why JSLR outperforms
811 the local structure learning method such as LPP and SpLpp
812 is that JSLR utilizes α -norm regularization for feature se-
813 lection and feature extraction to obtain the discriminative
814 information for face recognition.

- 815 3) As it can be seen from the Fig. 2(d), (e) and (f), the tra-
816 ditional regression methods and/or their extensions can
817 obtain only c projections for feature extraction and classi-
818 fication, which is not enough to achieve high recognition
819 rates. Note that the number of classes in CMU PIE and
820 Yale is 68 and 15, respectively. Therefore the numbers of
821 projections obtained by LDA are 67 and 14, and the num-
822 bers of projections obtained by RR are 68 and 15 on CMU
823 PIE and Yale databases, respectively. Fig. 2(e) and (f) show
824 that the recognition rates of RR and LDA achieve their top

TABLE XI
THE COMPUTATIONAL COST (UNIT: S) OF DIFFERENT METHODS

Data set (l)	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR
FERET ($l=4$)	0.0720	0.0390	14.7585	0.0201	1.7227	23.8032	0.1456	0.0642	1.7157	21.1985	0.0334	0.0301
AR ($l=4$)	0.0305	0.0566	7.2524	0.0097	0.5550	39.4835	0.1484	0.0607	1.1984	11.4715	0.0242	0.0237
CMU PIE ($l=4$)	0.0087	0.0767	7.6509	0.0375	0.1741	457.7060	0.0383	0.0483	0.4201	6.1862	0.0341	0.0906
Yale ($l=4$)	0.0019	0.0053	0.1162	0.0018	0.0161	0.4617	0.1251	0.0039	0.0505	1.1682	0.0028	0.0035
USPS ($l=400$)	0.0105	0.0433	4.0517	0.2033	131.6732	240.4646	0.0091	0.0376	6.7102	347.7577	0.1423	0.1462
PIE67 \times 170 ($l=80$)	0.1224	0.0738	168.2569	0.0289	285.8122	1.3378 $\times 10^3$	0.1190	0.2580	7.7944	630.7555	0.2123	0.2864
LFW ($l=4$)	0.0229	0.0661	48.6211	0.0157	1.4418	255.8663	0.0319	0.0984	1.1866	18.8623	0.1565	0.2012
Fa	0.0234	0.4945	1.3758 $\times 10^4$	0.0611	10.0917	3.8150 $\times 10^3$	0.0236	0.0371	18.1813	1.5484 $\times 10^3$	0.3679	0.2854

recognition rates using all the projections (we copy the final recognition rate to full fill all the dimensions listed on the horizontal axis). Thus the recognition rates no more increase after the number of dimension reaches 67 and 14 for LDA and 68 and 15 for RR on CMU PIE and Yale face databases, respectively. These figures show that the lack of enough projection of LDA and RR limits their performances. However, JSLR can break through this limitation and obtain more projections. This is the potential reason for JSLR to achieve higher recognition rates. In addition, the experimental result on USPS database with more than 4000 samples (as shown in Table VI) indicates the robustness and effectiveness of JSLR in dealing with large-sample size problem.

- 4) The $L_{2,1}$ -norm based methods such as JSLR, RFS and UDFS are robust to outliers in dataset and they guarantee the joint sparsity. However, JSLR obtains the best recognition rates when there are variations on lighting condition and face expressions. This indicates that JSLR is more robust than RFS and UDFS in feature extraction and selection when there exists variations on lighting condition and face expressions. In addition, the experimental results based on deep learning techniques presented in Table VIII and Table IX indicate the good performance of the proposed JSLR.
- 5) Experimental results indicate that the proposed JSLR performs better than the dictionary learning methods (LCLE-DL and LC-KSVD2). The reason is that JSLR guarantees the joint sparsity for discriminant feature selection or extraction in different cases. The Comparison between Deep-NN and Deep-JSLR shows that JSLR can further enhance the discriminative power of the deeply learned features based on CNNs in face recognition and character recognition tasks.
- 6) Table XI presents the computational time (unit: second) of each method on different data sets. From Table XI, we can know that the proposed JSLR based on $L_{2,1}$ -norm minimization is fast convergent and the computational cost is much less than the L_1 -norm based methods (i.e.

SLDA, SpLPP). The essential reason is that both SLDA and SpLPP use the least-angle regression method to compute the sparse solution and the iteration times are more than the proposed method. Moreover, the projections of SLDA and SpLPP are computed one by one while JSLR can simultaneously compute a set of jointly sparse projections. Thus, the proposed JSLR is efficient and effective for computer vision and pattern recognition.

- 7) From Tables III and IV, we can see that when more training samples are used, the recognition rates of all methods are higher than that when less training samples are used. However, it does not mean that more training samples can definitely help to obtain higher recognition rate. As shown in Table IV, when 20/24 training samples are used, the recognition rates of all methods become lower compared to the case when 19/24 training samples are used. The potential reason for this phenomenon is that too many training samples may lead to overfitting and thus all methods obtain poorer performance in the testing stage.

VII. CONCLUSION

Motivated by previous works that $L_{2,1}$ -norm regularization is able to obtain joint sparsity, and the local geometric information can enhance feature selection capability, in this paper, we propose a novel method called JSLR for feature extraction and selection. With $L_{2,1}$ -norm regularization and locality preserving property, JSLR can obtain any number of discriminative projections for feature selection, which addresses the drawback in LDA and ridge regression. Theoretical analyses show the close relationship of JSLR and ridge regression, which also guarantees the effectiveness of JSLR in feature extraction and selection. In order to obtain the optimal solution of JSLR, we propose an iterative algorithm which is proved to be convergent. In addition, the computational complexity of the algorithm is also presented. The performance of JSLR on several well-known face databases shows that it outperforms the classical principle component analysis methods, traditional sparse learning methods and recently proposed $L_{2,1}$ -norm regularization methods.

900 APPENDIX
 901 PROOF OF THEOREM 1

902 From Eq. (10), we have

$$903 \begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{Y} - 2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B}). \end{aligned}$$

903 By setting the derivatives of the above problem with respect
 904 to \mathbf{B} equaling to 0, we have

$$905 \mathbf{B} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}.$$

905 Let \mathbf{B}^* represents the optimal solution of Eq. (10), then

$$906 \mathbf{B}^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}.$$

906 Since $\mathbf{P}^0 = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}$, we have $\mathbf{B}^* = \mathbf{P}^0 \mathbf{A}$. As matrix \mathbf{A} is a rotation matrix, then the subspace spanned by \mathbf{B}^* in Eq. (10) is the same as that spanned by \mathbf{P}^0 in Eq. (1), namely, $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$.

907 Suppose $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, by the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, we have $\mathbf{B}^* = \mathbf{U} \frac{1}{\mathbf{D}} \mathbf{V}^T \mathbf{Y} \mathbf{A}$. Since the optimal solution of Eq. (3) is $\mathbf{P}^* = \mathbf{U} \frac{\mathbf{D}}{\mathbf{D}^2 + \alpha \mathbf{I}} \mathbf{V}^T \mathbf{Y}$, then we can find that the subspaces spanned by \mathbf{B}^* and \mathbf{P}^* have the same base matrix \mathbf{U} and the only difference is that there is a weighted rotation matrix, which does not affect the spanned subspace. Thus, we can say $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$.

907 If $\alpha \rightarrow 0$, we have

$$908 \mathbf{P}^* \mathbf{A} = \mathbf{U} \frac{\mathbf{D}}{\mathbf{D}^2 + \alpha \mathbf{I}} \mathbf{V}^T \mathbf{Y} \mathbf{A} \rightarrow \mathbf{U} \frac{1}{\mathbf{D}} \mathbf{V}^T \mathbf{Y} \mathbf{A} = \mathbf{B}^*.$$

908 Thus, for any two pattern vectors \mathbf{x}_i and \mathbf{x}_j , since $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, the distance of the two points obtained by using the two subspaces (i.e. \mathbf{B}^* and \mathbf{P}^*) for feature extraction is invariant. That is,

$$909 \|\mathbf{(x}_i - \mathbf{x}_j)^T \mathbf{B}^*\|_2 = \|\mathbf{(x}_i - \mathbf{x}_j)^T \mathbf{P}^* \mathbf{A}\|_2 = \|\mathbf{(x}_i - \mathbf{x}_j)^T \mathbf{P}^*\|_2,$$

910 which indicates that the performance of using the two metric matrices derived by \mathbf{B}^* and \mathbf{P}^* for classification will be the same.

910 For Eq. (19), if $\alpha_1 \rightarrow 0$ and $\alpha_2 \rightarrow 0$, we have

$$911 \begin{aligned} \bar{\mathbf{B}} &= (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} \\ &\rightarrow (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} = \mathbf{B}^*, \end{aligned}$$

912 namely, $\text{span}(\mathbf{B}^*) \rightarrow \text{span}(\bar{\mathbf{B}})$.

912 Since $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$, $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$, then $\text{span}(\bar{\mathbf{B}}) \rightarrow \text{span}(\mathbf{P}^0)$, $\text{span}(\bar{\mathbf{B}}) \rightarrow \text{span}(\mathbf{P}^*)$.

913 REFERENCES

- [1] G. V. Lashkia and L. Anthony, "Relevant, irredundant feature selection and noisy example elimination," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 34, no. 2, pp. 888–897, Apr. 2004.
- [2] T. W. S. Chow, P. Wang, and E. W. M. Ma, "A new feature selection scheme using a data distribution factor for unsupervised nominal data," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 38, no. 2, pp. 499–509, Apr. 2008.
- [3] J. Qian, J. Yang, and Y. Xu, "Local structure-based image decomposition for feature extraction with applications to face recognition," *IEEE Trans. Image Process.*, vol. 22, no. 9, pp. 3591–3603, Sep. 2013.
- [4] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *J. Comput. Graph. Stat.*, vol. 15, no. 2, pp. 265–286, 2006.
- [5] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. fisherfaces: Recognition using class specific linear projection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 7, pp. 711–720, Jul. 1997.
- [6] W. Zheng and X. Tang, "A robust algorithm for generalized orthonormal discriminant vectors," in *Proc. IEEE 18th Int. Conf. Pattern Recognit.*, 2006, vol. 2, pp. 784–787.
- [7] G. Dai and Y. Qian, "A gabor direct fractional-step LDA algorithm for face recognition," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2004, vol. 1, pp. 61–64.
- [8] I. Dagher, "Incremental PCA-LDA algorithm," in *Proc. IEEE Int. Conf. Comput. Intell. Meas. Syst. Appl.*, 2010, pp. 97–101.
- [9] R. Raghavendra, B. Dorizzi, A. Rao, and G. H. Kumar, "Designing efficient fusion schemes for multimodal biometric systems using face and palmprint," *Pattern Recognit.*, vol. 44, no. 5, pp. 1076–1088, 2011.
- [10] M. Eskandari and Ö. Toygar, "Selection of optimized features and weights on face-Iris fusion using distance images," *Comput. Vision Image Understanding*, vol. 137, pp. 63–75, 2015.
- [11] F. Zhang, J. Yang, J. Qian, and Y. Xu, "Nuclear norm-based 2-DPCA for extracting features from images," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 10, pp. 2247–2260, Oct. 2015.
- [12] W. Wang, Y. Yan, F. Nie, S. Yan, and N. Sebe, "Flexible manifold learning with optimal graph for image and video representation," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2664–2675, Jun. 2018.
- [13] F. Nie, S. Yang, R. Zhang, and X. Li, "A general framework for auto-weighted feature selection via global redundancy minimization," *IEEE Trans. Image Process.*, vol. 28, no. 5, pp. 2428–2438, Dec. 2018.
- [14] H. Zhang, J. Yang, J. Xie, J. Qian, and B. Zhang, "Weighted sparse coding regularized nonconvex matrix regression for robust face recognition," *Inf. Sci.*, vol. 394, pp. 1–17, 2017.
- [15] X. Ning, W. Li, B. Tang, and H. He, "BULDP: Biomimetic uncorrelated locality discriminant projection for feature extraction in face recognition," *IEEE Trans. Image Process.*, vol. 27, no. 5, pp. 2575–2586, Feb. 2018.
- [16] R. Tibshirani, "Regression shrinkage and selection via the lasso: A retrospective," *J. Roy. Stat. Soc.: Ser. B (Stat. Method.)*, vol. 73, no. 3, pp. 273–282, 2011.
- [17] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. Roy. Stat. Soc.: Ser. B (Stat. Method.)*, vol. 67, no. 2, pp. 301–320, 2005.
- [18] I. T. Jolliffe, N. T. Trendafilov, and M. Uddin, "A modified principal component technique based on the lasso," *J. Comput. Graph. Stat.*, vol. 12, no. 3, pp. 531–547, 2003.
- [19] A. d'Aspremont, L. E. Ghaoui, M. I. Jordan, and G. R. Lanckriet, "A direct formulation for sparse PCA using semidefinite programming," in *Proc. Adv. Neural Inf. Process. Syst.*, 2005, pp. 41–48.
- [20] J. Feng, L. Jiao, F. Liu, T. Sun, and X. Zhang, "Unsupervised feature selection based on maximum information and minimum redundancy for hyperspectral images," *Pattern Recognit.*, vol. 51, pp. 295–309, 2016.
- [21] Z. Li and J. Tang, "Unsupervised feature selection via nonnegative spectral analysis and redundancy control," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5343–5355, Dec. 2015.
- [22] L. Clemmensen, T. Hastie, D. Witten, and B. Ersbøll, "Sparse discriminant analysis," *Technometrics*, vol. 53, no. 4, pp. 406–413, 2011.
- [23] Z. Qiao, L. Zhou, and J. Z. Huang, "Sparse linear discriminant analysis with applications to high dimensional low sample size data," *Int. J. Appl. Math.*, vol. 39, no. 1, 2009.
- [24] Z. Zeng, X. Wang, J. Zhang, and Q. Wu, "Semi-supervised feature selection based on local discriminative information," *Neurocomputing*, vol. 173, pp. 102–109, 2016.
- [25] C. Shi, Q. Ruan, G. An, and R. Zhao, "Hessian semi-supervised sparse feature selection based on $l_{2,1/2}$ -matrix norm," *IEEE Trans. Multimedia*, vol. 17, no. 1, pp. 16–28, Jan. 2014.
- [26] C. Ding and D. Tao, "Robust face recognition via multimodal deep face representation," *IEEE Trans. Multimedia*, vol. 17, no. 11, pp. 2049–2058, Nov. 2015.
- [27] A. Majumdar and R. K. Ward, "Classification via group sparsity promoting regularization," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, 2009, pp. 861–864.
- [28] A. Majumdar and R. K. Ward, "Robust classifiers for data reduced via random projections," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 40, no. 5, pp. 1359–1371, Oct. 2010.
- [29] A. Majumdar and R. K. Ward, "Fast group sparse classification," *Can. J. Electr. Comput. Eng.*, vol. 34, no. 4, pp. 136–144, 2009.

- 1015 [30] X. He, D. Cai, S. Yan, and H.-J. Zhang, "Neighborhood preserving embedding," in *Proc. 10th IEEE Int. Conf. Comput. Vision Volume 1*, 2005, vol. 2, pp. 1208–1213.
- 1016 [31] D. Cai *et al.*, "Isometric projection," in *Proc. Assoc. Advancement Artif. Intell.*, 2007, pp. 528–533.
- 1017 [32] X. He and P. Niyogi, "Locality preserving projections," in *Proc. Advances Neural Inf. Process. Syst.*, 2004, pp. 153–160.
- 1018 [33] X. He, S. Yan, Y. Hu, P. Niyogi, and H.-J. Zhang, "Face recognition using laplacianfaces," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 3, pp. 328–340, Mar. 2005.
- 1019 [34] D. Cai, X. He, and J. Han, "Spectral regression: A unified approach for sparse subspace learning," in *Proc. IEEE Int. Conf. Data Mining*, 2007, pp. 73–82.
- 1020 [35] X. Liu, L. Wang, J. Zhang, J. Yin, and H. Liu, "Global and local structure preservation for feature selection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 6, pp. 1083–1095, Jun. 2014.
- 1021 [36] S. Liao *et al.*, "Discriminant analysis via joint euler transform and $\ell_{2,1}$ -norm," *IEEE Trans. Image Process.*, vol. 27, no. 11, pp. 5668–5682, Nov. 2018.
- 1022 [37] J. Yang *et al.*, "Nuclear norm based matrix regression with applications to face recognition with occlusion and illumination changes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 1, pp. 156–171, Jan. 2017.
- 1023 [38] J. Qian, L. Lei, Y. Jian, F. Zhang, and Z. Lin, "Robust nuclear norm regularized regression for face recognition with occlusion," *Pattern Recognit.*, vol. 48, no. 10, pp. 3145–3159, 2015.
- 1024 [39] L. Luo, J. Yang, J. Qian, Y. Tai, and G. F. Lu, "Robust image regression based on the extended matrix variate power exponential distribution of dependent noise," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 9, pp. 2168–2182, Sep. 2017.
- 1025 [40] L. Luo, J. Yang, J. Qian, and Y. Tai, "Nuclear- ℓ_1 norm joint regression for face reconstruction and recognition with mixed noise," *Pattern Recognit.*, vol. 48, no. 12, pp. 3811–3824, 2015.
- 1026 [41] J. Chen, J. Yang, L. Luo, J. Qian, and W. Xu, "Matrix variate distribution-induced sparse representation for robust image classification," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 10, pp. 2291–2300, Oct. 2015.
- 1027 [42] M. Yang, L. Zhang, J. Yang, and D. Zhang, "Regularized robust coding for face recognition," *IEEE Trans. Image Process.*, vol. 22, no. 5, pp. 1753–1766, May 2013.
- 1028 [43] Y. Li, J. Si, G. Zhou, S. Huang, and S. Chen, "FREL: A stable feature selection algorithm," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 7, pp. 1388–1402, Jul. 2015.
- 1029 [44] W. Yang, C. Sun, and W. Zheng, "A regularized least square based discriminative projections for feature extraction," *Neurocomputing*, vol. 175, pp. 198–205, 2016.
- 1030 [45] H. Pu and G. Gao, "Parameterless reconstructive discriminant analysis for feature extraction," *Neurocomputing*, vol. 190, pp. 50–59, 2016.
- 1031 [46] Y. Yang, Z.-J. Zha, Y. Gao, X. Zhu, and T.-S. Chua, "Exploiting web images for semantic video indexing via robust sample-specific loss," *IEEE Trans. Multimedia*, vol. 16, no. 6, pp. 1677–1689, Oct. 2014.
- 1032 [47] S. Xiang, F. Nie, G. Meng, C. Pan, and C. Zhang, "Discriminative least squares regression for multiclass classification and feature selection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 23, no. 11, pp. 1738–1754, Nov. 2012.
- 1033 [48] J. Han, Z. Sun, and H. Hao, " ℓ_0 -norm based structural sparse least square regression for feature selection," *Pattern Recognit.*, vol. 48, no. 12, pp. 3927–3940, 2015.
- 1034 [49] J. Yang and C. Ong, "An effective feature selection method via mutual information estimation," *IEEE Trans. Syst. Man Cybern. Part B Cybern.*, vol. 42, no. 6, pp. 1550–1559, Dec. 2012.
- 1035 [50] Z. Zeng, H. Zhang, R. Zhang, and C. Yin, "A novel feature selection method considering feature interaction," *Pattern Recognit.*, vol. 48, no. 8, pp. 2656–2666, 2015.
- 1036 [51] Y. Hui and Y. Jian, "Sparse discriminative feature selection," *Pattern Recognit.*, vol. 48, no. 5, pp. 1827–1835, 2015.
- 1037 [52] Y. Liu, F. Tang, and Z. Zeng, "Feature selection based on dependency margin," *IEEE Trans. Cybern.*, vol. 45, no. 6, pp. 1209–1221, Jun. 2015.
- 1038 [53] C. Hou, F. Nie, X. Li, D. Yi, and Y. Wu, "Joint embedding learning and sparse regression: A framework for unsupervised feature selection," *IEEE Trans. Cybern.*, vol. 44, no. 6, pp. 793–804, Jun. 2014.
- 1039 [54] F. Nie, H. Huang, C. Xiao, and C. H. Q. Ding, "Efficient and robust feature selection via joint $\ell_2, 1$ -norms minimization," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2010, pp. 1813–1821.
- 1040 [55] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*. Berlin/Heidelberg: Springer, Jul. 2003.
- 1041 [56] Y. Yang, H. T. Shen, Z. Ma, Z. Huang, and X. Zhou, "L_{2,1}-norm regularized discriminative feature selection for unsupervised learning," in *Proc. Int. Joint Conf. Artif. Intell.*, 2011, pp. 1589–1594.
- 1042 [57] Y. Yang, Z. Ma, Y. Yang, F. Nie, and H. T. Shen, "Multitask spectral clustering by exploring intertask correlation," *IEEE Trans. Cybern.*, vol. 45, no. 5, pp. 1083–1094, May 2015.
- 1043 [58] E. V. Den Berg and M. P. Friedlander, "Theoretical and empirical results for recovery from multiple measurements," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2516–2527, May 2010.
- 1044 [59] A. Majumdar and R. K. Ward, "Synthesis and analysis prior algorithms for joint-sparse recovery," in *Proc. IEEE Int. Conf. Acoust.*, 2012, pp. 3421–3424.
- 1045 [60] J. He, L. Ding, L. Jiang, and L. Ma, "Kernel ridge regression classification," in *Proc. Int. Joint Conf. Neural Netw.*, 2014, pp. 2263–2267.
- 1046 [61] Z. Zheng *et al.*, "Regression analysis of locality preserving projections via sparse penalty," *Inf. Sci.*, vol. 303, pp. 1–14, 2015.
- 1047 [62] G. Shikkenawis and S. K. Mitra, "On some variants of locality preserving projection," *Neurocomputing*, vol. 173, pp. 196–211, 2016.
- 1048 [63] F. Nie, X. Wang, M. I. Jordan, and H. Huang, "The constrained laplacian rank algorithm for graph-based clustering," in *Proc. 30th AAAI Conf. Artif. Intell.*, 2016, pp. 1969–1976.
- 1049 [64] H. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos, "Multilinear principal component analysis of tensor objects for recognition," in *Proc. Int. Conf. Pattern Recognit.*, 2006, pp. 776–779.
- 1050 [65] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vision*, 2014, pp. 818–833.
- 1051 [66] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv preprint arXiv:1409.1556*.
- 1052 [67] Q. Wang, Z. Qin, F. Nie, and Y. Yuan, "Convolutional 2D LDA for nonlinear dimensionality reduction," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 2929–2935.
- 1053 [68] Z. Jiang, L. Zhe, and L. S. Davis, "Learning a discriminative dictionary for sparse coding via label consistent K-SVD," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2011, pp. 1697–1704.
- 1054 [69] Z. Li, Z. Lai, Y. Xu, J. Yang, and D. Zhang, "A locality-constrained and label embedding dictionary learning algorithm for image classification," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 2, pp. 278–293, Feb. 2017.
- 1055 [70] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *Proc. Eur. Conf. Comput. Vision*, 2016, pp. 499–515.
- 1056 [71] P. J. Phillips, H. Moon, S. A. Rizvi, and P. J. Rauss, "The FERET evaluation methodology for face-recognition algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 10, pp. 1090–1104, Oct. 2000.
- 1057 [72] A. Martinez and R. Benavente, "The AR face database," CVC, West Lafayette, IN, USA, Tech. Rep. 24, 1998.
- 1058 [73] T. Sim, S. Baker, and M. Bsat, "The CMU pose, illumination and expression database of human faces," Carnegie Mellon Univ.: Pittsburgh, Pennsylvania, USA, Tech. Rep. CMU-RI-TR-OI-02, 2001.
- 1059 [74] G. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," Tech. Rep. 07-49, UMass, 2007.
- 1060 [75] Y. Jia *et al.*, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. 22nd ACM Int. Conf. Multimedia*, 2014, pp. 675–678.

Jointly Sparse Locality Regression for Image Feature Extraction

Dongmei Mo^{ID}, Zhihui Lai^{ID}, Xizhao Wang^{ID}, *Fellow, IEEE*, and Waikeung Wong^{ID}

Abstract—This paper proposes a novel method called Jointly Sparse Locality Regression (JSLR) for feature extraction and selection. JSLR utilizes joint $L_{2,1}$ -norm minimization on regularization term, and also introduces the locality to characterize the local geometric structure of the data. There are three main contributions in JSLR for face recognition. Firstly, it eliminates the drawback in ridge regression and Linear Discriminant Analysis (LDA) that when the number of the classes is too small, not enough projections can be obtained for feature extraction. Secondly, by using the local geometric structure as the regularization term, JSLR is able to preserve local information and find an embedding subspace which can detect the most essential data manifold structure. Moreover, since the $L_{2,1}$ -norm based loss function is robust to outliers in data points, JSLR provides the joint sparsity for robust feature selection. The theoretical connections of the proposed method and the previous regression methods are explored and the convergence of the proposed algorithm is also proved. Experimental evaluation on several well-known data sets shows the merits of the proposed method on feature selection and classification.

Index Terms—Regression, face recognition, feature extraction, local structure, joint sparsity.

I. INTRODUCTION

SINCE the data used in computer vision or pattern recognition is very high dimensional, it is of great importance to select the key features from large quantities of variables. Besides, the redundancy of the data would affect the performance of some algorithms in practical applications [1], and thus most of the algorithms cannot obtain a good performance in high-dimensional case [2]. Therefore, feature extraction and selection are of great importance in processing the high-dimensional data set [3].

Manuscript received May 4, 2019; revised September 9, 2019 and December 11, 2019; accepted December 13, 2019. This work was supported in part by The Hong Kong Polytechnic University (Project Code: RHR1) and in part by General Research Fund of the Research Grants Council of Hong Kong (Project Code: 15202217). The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jingdong Wang. (*Corresponding author: Waikeung Wong*.)

D. Mo and W. Wong are with the Institute of Textiles and Clothing, The Hong Kong Polytechnic University, Hong Kong SAR of China (e-mail: dongmei.mo@connect.polyu.hk; calvin.wong@polyu.edu.hk).

Z. Lai is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China, and also with the Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen 518060, China (e-mail: lai_zhi_hui@163.com).

X. Wang is with the College of Computer Science and Software Engineering and Guangdong Key Lab of Intelligent Information Processing, Shenzhen University, Shenzhen 518060, China (e-mail: xizhaowang@ieee.org).

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2019.2961508

Up to now, one of the classical methods is Principle Component Analysis (PCA) [4], which is a simple and effective unsupervised method as it solves the eigen decomposition problem to obtain the optimal vectors for dimensionality reduction. Linear Discriminant Analysis (LDA) [5] is a representative supervised method in feature extraction and dimensionality reduction, which uses the label information to improve the performance in classification. By maximizing the ratio of the between-class scatter to the within-class scatter of the training dataset, LDA can obtain an optimal set of discriminative vectors [6]. However, a drawback of LDA is that it is unsuitable for small sample size problem in face recognition. An effective model called PCA+LDA [7], which joints the two major techniques to obtain the discriminant vectors [8], has been proposed to deal with the problem. The other two methods called Particle Swarm Optimization (PSO) [9] and Backtracking Search Algorithm (BSA) [10] can also significantly reduce the number of features so as to reduce the computational complexity and at the same time guarantee the same level of performance.

However, PCA and LDA cannot provide the sparse projections for feature extraction since the learned projections are the linear combination of the data [11]. Recently, sparse regression showed the outstanding performance in feature selection and extraction [12]–[15]. By adding sparsity penalty for feature selection, the accuracy and robustness of these methods might be improved. Thus many studies focused on the sparse learning for variable selection. Zou *et al.* proposed an effective model called Sparse Principle Component Analysis (SPCA) [4] to generate modified principle components with sparse loadings by using the lasso or elastic net constraint [16], [17]. Some other sparse PCA algorithms, such as the SCoTLASS algorithm [18], the DSPCA algorithm [19] were proposed. All of these methods focused on sparse learning without using the class label information. Besides, Feng *et al.* proposed the unsupervised learning method based on maximum information and minimum redundancy (MIMR) [20] for hyperspectral image analysis, and Li *et al.* proposed another unsupervised feature selection method by nonnegative spectral analysis and redundancy control [21].

Some other researchers developed the supervised methods using the label information to perform sparse learning for feature extraction and selection. One of the effective methods is Sparse Discriminant Analysis (SDA) [22], which extends linear discriminant analysis to sparse case by imposing the sparsity constraint. Moreover, to overcome the data piling problem of LDA in the high dimensional and low sample size (HDLSS) case, Qiao *et al.* proposed sparse LDA to obtain sparse linear

discriminant vectors by taking the relationship between Fisher's LDA and a generalized eigenvalue problem into consideration [23]. Besides, some semi-supervised methods were also proposed. A semi-supervised method, which used partially labeled data samples, was designed to achieve batch feature selection [24]. Another method called Hessian sparse feature selection based on $L_{2,1/2}$ -matrix norm (HFSL) was proposed for semi-supervised sparse feature selection [25]. For the multimodal case, Ding *et al.* proposed a method using multimodal information to jointly learn face representation [26]. Furthermore, sparse regularization learning were also used in classification designation for different pattern recognition tasks [27]–[29].

It is a well-known fact that not all data are distributed on a linear subspace. They may lie on the nonlinear low-dimensional manifold embedding on the high-dimensional ambient space. Therefore, many manifold learning algorithms were proposed. The representative methods include Neighborhood Preserving Embedding (NPE) [30], Isometric Projection (IsoP) [31] and Locality Preserving Projection(LPP) [32], [33], etc. These algorithms aimed to preserve the local geometric structure of the data manifold. By introducing the locality for sparse subspace learning, Cai *et al.* also proposed a new method called Unified Sparse Subspace Learning (USSL) [34]. USSL utilized the elastic net for regression to simultaneously select the most important variables and take the local geometric structure into consideration. Besides, by combining the global pairwise sample similarity with local geometric structure, a new method called GLSPFS [35] was proposed by Liu *et al.* for feature selection.

In recent years, a great deal of attention has been paid to the regression methods with different norms for image recognition, feature extraction and variable selection [36]. For example, nuclear norm regression methods were proposed in [37], [38] for face recognition. The L_1 -norm based sparse regularized learning methods [39]–[41] have been used for face reconstruction and recognition [42]. A feature selection algorithm framework called Feature-weighting as Regularized Energy-based Learning (FREL) was proposed by Li *et al.* [43]. Based on least square regularization, Yang *et al.* [44] proposed the discriminative projection method. And the traditional RDA was further developed as Parameterless Reconstructive Discriminant Analysis (PRDA) [45] for feature extraction. In [46], the L_1 -norm minimization was employed to design a specific loss function, by which the abundant user tagged Web images are treated as noisy samples and will not be emphasized so as to perform robust semantic video indexing. Other methods, such as [47]–[53] were also proposed to deal with different feature selection problems. The methods in terms of jointly sparse subspace learning attracted great attention in the field of feature selection. Since the $L_{2,1}$ -norm based regression loss function is robust to outlier in data set, it can improve the robustness in learning steps. Therefore, some algorithms with joint $L_{2,1}$ -norm regularization were proposed to guarantee the joint sparsity for feature extraction. The model called Robust Feature Selection (RFS) [54] via joint $L_{2,1}$ -norms minimization showed the good performance for feature selection with joint sparsity. Yang *et al.* proposed another model called Unsupervised Discriminative Feature Selection (UDFS) [55] for sparse subspace learning. Experimental

results showed that UDFS outperforms the existing unsupervised feature extraction methods and its main advantage is that UDFS not only uses discriminative information but also uses local structure of datas distribution for feature selection [56]. The $L_{2,1}$ -norm regularization is also used in [57] to discover the common features shared across all the clustering tasks so as to obtain a discriminative low dimensional space for clustering. Except for the jointly sparse feature selection, the $L_{2,1}$ -norm was also widely used to deal with the joint-sparse recovery problems [58], [59] in computer vision.

Although a lot of methods have been developed to improve the performance of regression methods, there still exist some problems to be solved. For example, when the number of the class is too small, not enough projections can be obtained by the classical regression methods and/or their extensions to achieve higher classification accuracy. Also, most existing regression method do not simultaneously consider the geometric structure of the data as well as the sparsity of the projection matrix. In this paper, we propose a novel model called Jointly Sparse Locality Regression (JSLR) for feature extraction and selection. JSLR can not only avoid the limitation in the existing regression methods but also guarantee the sparsity by using $L_{2,1}$ -norm regularization on the projection matrix. What is more, JSLR incorporates the local structure of the data in regression form, by which the optimization problem can be easily optimized so as to obtain better performance of feature extraction with less computational time.

The main contributions of this paper are described as below:

- 1) The number of the projections in LDA-based methods or regression-based methods is limited by the rank of the so-called between-class scatter matrix or the number of the classes. The proposed method can break out the limitation to obtain more projections for feature extraction by designing a novel regression model.
- 2) Theoretical connections between the proposed method and the previous regression methods are discovered. Moreover, the convergence of the proposed algorithm is also proved.
- 3) The experimental results of the proposed model with or without $L_{2,1}$ -norm regularization indicate that adding $L_{2,1}$ -norm penalty on the projection matrix can obtain joint sparsity for feature extraction so as to achieve high recognition rate.

The rest of this paper is organized as follows: In Section II, we discuss the related works and the extension based on ridge regression will be shown in Section III. In Section IV, we propose our objective function and the local optimal solution. Section V focuses on theoretical analysis (the convergence and the computational complexity). The proposed model will be evaluated by several well-known databases in Section VI. In Section VII, we draw a conclusion for this paper.

II. RELATED WORKS

In this section, the notations used in this paper will be briefly described and the related works will be reviewed.

192 *A. Notations*

193 Scalars are denoted as lowercase or uppercase italic letters,
 194 i.e. $i, j, d, p, n, c, \alpha_1, \alpha_2$ etc. while vectors are represented as
 195 bold lowercase italic letters, i.e. \mathbf{x}, \mathbf{y} , etc. Matrices are defined
 196 as bold uppercase italic letters, i.e. $\mathbf{A}, \mathbf{B}, \mathbf{X}, \mathbf{Y}, \mathbf{W}$ etc.

197 Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in R^{d \times n}$ then \mathbf{X} denotes a $d \times n$
 198 matrix as the original data set, where n is the number of total
 199 training samples and d denotes the features dimension for each
 200 sample. Let $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_c] \in R^{n \times c}$ to be a $n \times c$ matrix
 201 as the label matrix falling into c classes.

202 *B. Regressions*

203 Ridge Regression [60] is a regularized least square method
 204 for multivariate learning. It aims to solve the multicollinearity
 205 problem of covariates in samples.

206 The optimization problem of the simplest regression is

$$\mathbf{P}^0 = \arg \min_{\mathbf{P}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2 \quad (1)$$

207 where \mathbf{X} denotes the training set of n training data. The ma-
 208 trix $\mathbf{P} \in R^{d \times c}$ aims to lead the linear dependency between the
 209 training data and the corresponding labels. By setting the deriva-
 210 tives of (1) with respect to \mathbf{P} equaling to 0, we have the optimal
 211 solution

$$\mathbf{P}^0 = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \quad (2)$$

212 However this optimal solution is only suitable for the case when
 213 $\mathbf{X} \mathbf{X}^T$ is a full-rank matrix. Because of the small-sample size
 214 problem, the matrix $\mathbf{X} \mathbf{X}^T$ may be not a full-rank one. Therefore,
 215 to solve the singular problem in computing the inverse of
 216 $\mathbf{X} \mathbf{X}^T$ the L_2 -norm regularized term was added to (1), and then
 217 we have the classical ridge regression optimization problem:

$$\mathbf{P}^* = \arg \min_{\mathbf{P}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2 + \alpha \|\mathbf{P}\|_F^2 \quad (3)$$

218 By setting the derivatives of (3) with respect to \mathbf{P} equaling to 0,
 219 we have the optimal solution for (3) as

$$\mathbf{P}^* = (\mathbf{X} \mathbf{X}^T + \alpha \mathbf{I})^{-1} \mathbf{X} \mathbf{Y} \quad (4)$$

220 For further analysis in the following sections, we need to rep-
 221 resent the optimal solution of (3). Based on the SVD of $\mathbf{X} =$
 222 $\mathbf{U} \mathbf{D} \mathbf{V}^T$, the optimal solution can be represented as

$$\mathbf{P}^* = \mathbf{U} \frac{\mathbf{D}}{\mathbf{D}^2 + \alpha \mathbf{I}} \mathbf{V}^T \mathbf{Y} \quad (5)$$

223 From (2) and (5), we can know that the optimal projection matrix
 224 \mathbf{P}^0 and \mathbf{P}^* have the size, i.e. $d \times c$. That is, we can obtain only
 225 c projective vectors for feature extraction.

226 *C. The Review of LPP*

227 LPP [32], [33] computes the best linear approximations to
 228 the eigenfunctions of the manifold's Laplace Beltrami opera-
 229 tor. It aims to preserve local information and to find an embed-
 230 ding subspace which detects the most essential data manifold

231 structure [61], [62]. The objective function of LPP is to mini-
 232 mize

$$\begin{aligned} \frac{1}{2} \sum_{ij} (\mathbf{y}_i - \mathbf{y}_j)^2 \bar{w}_{ij} &= \frac{1}{2} \sum_{ij} \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|^2 \bar{w}_{ij} \\ &= \text{tr} (\mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B}) = \text{tr} (\mathbf{B}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{B}) \end{aligned} \quad (6)$$

233 where matrix \mathbf{B} denotes the transformation matrix, \mathbf{y}_i and \mathbf{y}_j
 234 denote the low-dimensional vectors of \mathbf{x}_i and \mathbf{x}_j in subspace
 235 \mathbf{B} , respectively. $\bar{\mathbf{W}}$ is supposed to be the similarity matrix of
 236 all pairwise data points, $\mathbf{L} = \bar{\mathbf{D}} - \bar{\mathbf{W}}$ is Laplacian matrix. $\bar{\mathbf{D}}$
 237 is a diagonal matrix and its element \bar{d}_{ii} is column or row sum
 238 of matrix $\bar{\mathbf{W}}$ (because $\bar{\mathbf{W}}$ is symmetric), i.e. $\bar{d}_{ii} = \sum_i \bar{w}_{ij}$.
 239

The similarity matrix \bar{w}_{ij} is defined as:

$$\bar{w}_{ij} = \begin{cases} \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/t), & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\|^2 < \varepsilon, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

240 where parameter $t \in R$, ε denotes the radius of the local neigh-
 241 borhood and it can be a sufficiently small positive value ($\varepsilon > 0$).
 242 In Eq. (7), the similarity matrix \bar{w}_{ij} might be sensitive to the
 243 value of the parameter t . To solve this problem, recently a
 244 parameter-free method was proposed in [63].
 245

246 By considering the similarity matrix $\bar{\mathbf{W}}$, the relationship be-
 247 tween each data pair \mathbf{x}_i and \mathbf{x}_j in original space can be preserved
 248 by reconstructing the relationship between \mathbf{y}_i and \mathbf{y}_j in the low
 249 dimensional space \mathbf{B} with $\sum_{ij} (\mathbf{y}_i - \mathbf{y}_j)^2 \bar{w}_{ij}$ where $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i$
 250 and $\mathbf{y}_j = \mathbf{B}^T \mathbf{x}_j$. The optimal projections can be obtained by
 251 solving the following generalized eigen-function:

$$\mathbf{X} \mathbf{L} \mathbf{X}^T b = \lambda \mathbf{X} \bar{\mathbf{D}} \mathbf{X}^T. \quad (8)$$

252 Suppose $\lambda_i (i = 1, 2, \dots, d)$ are eigenvalues of problem 8,
 253 we can sort the eigenvalues in ascending order, then matrix
 254 $\mathbf{B} = [\mathbf{b}^1, \mathbf{b}^2, \dots, \mathbf{b}^k]$ combined of k eigenvectors corresponding
 255 to the first k smallest eigenvalues is the final projection matrix
 of LPP.

256

III. THE EXTENSION BASED ON RIDGE REGRESSION

257 In this section, we firstly review the definition of $L_{2,1}$ -norm
 258 and its property. Then we analyze the advantages and disadvan-
 259 tages of ridge regression. Meanwhile, we also propose a simple
 260 extension based on ridge regression.

261 *A. The Definition of $L_{2,1}$ -Norm and Its Property*

262 Some well-known models such as PCA, multilinear PCA
 263 (MPCA) [64], etc. use L_2 -norm as the measurement to com-
 264 pute the optimal projections in computer vision and face recog-
 265 nition. However, a large amount of experimental results have
 266 shown that in sparse feature selection, L_1 -norm outperforms
 267 L_2 -norm because of its generalization and the robustness for
 268 classification [16], [17], [61]. By combining the advantages of
 269 both L_1 -norm and part property of L_2 -norm, researchers obtain
 270 joint $L_{2,1}$ -norm minimization on both loss functions and regu-
 271 larization term for robust sparse learning for feature extraction
 272 [53]. Therefore, we use the $L_{2,1}$ -norm instead of L_2 -norm as a

273 new measurement for model design to overcome the problem of
 274 L_2 -norm being sensitive to outliers in a certain sense [54].

275 The $L_{2,1}$ -norm of a matrix is defined as

$$\|\mathbf{M}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^m m_{ij}^2} = \sum_{i=1}^n \|\mathbf{m}^i\|_2 \quad (9)$$

276 where the i -th row and the j -th column of a matrix $\mathbf{M} = (\mathbf{m}_i)_j$
 277 are denoted as \mathbf{m}^i and \mathbf{m}_j .

278 The common advantage of $L_{2,1}$ -norm and L_1 -norm based
 279 loss function is that they are more robust to outliers. However,
 280 the major difference between $L_{2,1}$ -norm and L_1 -norm is that
 281 $L_{2,1}$ -norm regularization is suitable for selecting meaningful or
 282 more powerful discriminant features from the data points with
 283 joint sparsity. The $L_{2,1}$ -norm based regularized methods can
 284 eliminate those useless interferences via making the elements
 285 in some rows of the projection matrix become zero such that
 286 the important features of the data points are emphasized and the
 287 insignificant features are ignored (filtered out) when conducting
 288 feature selection or extraction. Another advantage of $L_{2,1}$ -norm
 289 is that the $L_{2,1}$ -norm based methods are fast convergent and
 290 thus the computational cost is lower (this can be verified from
 291 computational cost of the $L_{2,1}$ -norm based methods compared
 292 with the L_1 -norm based methods in Table XI in Experiment
 293 section) [54], [56].

294 In all, employing the $L_{2,1}$ -norm instead of L_1 -norm as the
 295 regularization can obtain the joint sparsity to improve the per-
 296 formance and at the same time reduce the computational cost for
 297 efficient feature extraction and selection on image recognition
 298 tasks [54].

299 B. A Key Drawback in Traditional Regression

300 In (1) and (3), there exists a problem that when the number
 301 of the classes is too small, the traditional models cannot obtain
 302 enough projections for achieving good performance in pattern
 303 recognition. Thus, it is possible that learning more projection
 304 may improve the performance in feature extraction and classi-
 305 fication [4]. In order to obtain more projections in the regres-
 306 sion model, a tractable approach is to modify the representation
 307 $\|\mathbf{Y} - \mathbf{X}^T \mathbf{P}\|_F^2$ to be $\|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2$, which means that the
 308 matrix $(\mathbf{B} \mathbf{A}^T) \in R^{d \times c}$ takes the place of the matrix $\mathbf{P} \in R^{d \times c}$
 309 in the model. Thus we have the following optimization problem:

$$(\mathbf{A}^*, \mathbf{B}^*) = \arg \min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2, \text{ s.t. } \mathbf{A}^T \mathbf{A} = \mathbf{I}, \quad (10)$$

310 where \mathbf{A} is a $c \times k$ matrix and the size of matrix \mathbf{B} is $d \times k$
 311 where the notation k is any positive integer and c denotes the
 312 number of classes. In other words, the optimal solution \mathbf{B} with
 313 size $d \times k$ is able to break out the limitation of class number of
 314 the training data since the size of \mathbf{B} is not related to the class
 315 number and the variable k in \mathbf{B} is not related to the class number
 316 and the variable k in \mathbf{B} can be set as value that is larger than
 317 c , while \mathbf{P} with size $d \times c$ indicates that it can obtain at most c
 318 projections for feature selection.

From (10), we have

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{Y}) - \text{Tr}(2 \mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} - \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B}). \end{aligned} \quad (11)$$

319 By setting the derivatives of (11) with respect to \mathbf{B} equaling to
 320 0, the problem (11) is minimized at

$$\mathbf{B}^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}, \quad (12)$$

321 where \mathbf{B}^* represents the optimal solution of (10).

322 Denote the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, where $\mathbf{A}^T \mathbf{A} = \mathbf{I}$,
 323 $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, we have

$$\mathbf{B}^* = \mathbf{U} \frac{1}{\mathbf{D}} \mathbf{V}^T \mathbf{Y} \mathbf{A}. \quad (13)$$

324 For (1) and 10), we have following propositions:

325 *Proposition 1:* Suppose $\mathbf{X} \mathbf{X}^T$ is the full-rank matrix. Let
 326 \mathbf{P}^0 be the optimal solution to (1) and \mathbf{B}^* be the optimal solution
 327 to 10), if $k = c$ (i.e. the number of projection is equal to the
 328 number of class), then $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$.

329 If $k > c$, the optimal solution \mathbf{B}^* with size $d \times k$ in (10) can
 330 obtain k projections instead of c projections as obtained by \mathbf{P}^0 in
 331 (1), which breaks out the small-class problem. In Proposition 1,
 332 the reason why the small-class problem is addressed by (10) is
 333 that the \mathbf{P}^0 in (1) has c projections while the optimal solution
 334 \mathbf{B}^* for (10) can learn k projections to perform feature extraction
 335 and classification, where k can be set as any integer. In other
 336 words, the number of the learned projections from (10) is not
 337 limited by the number of class and thus the small-class problem
 338 is addressed.

339 Similarly, for (3) and 10), we have following proposition:

340 *Proposition 2:* Let \mathbf{P}^* be the optimal solution to (3) and \mathbf{B}^*
 341 be the optimal solution to 10), if $k = c$, (i.e. the number of projec-
 342 tion is equal to the number of class), then $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$.
 343 Furthermore if $\alpha \rightarrow 0$, the metric matrices derived by \mathbf{B}^* and
 344 \mathbf{P}^* for classification are equivalent to each other.

345 If $k > c$, the optimal solution \mathbf{B}^* with size $d \times k$ in (10) can
 346 obtain k projections instead of c projections as obtained by \mathbf{P}^*
 347 in (3), which breaks out the small-class problem.

348 Proposition 2 indicates that when $\alpha \rightarrow 0$ (or using $\alpha = \varepsilon$,
 349 where ε is a very small number), the performance using \mathbf{B}^* and
 350 \mathbf{P}^* for feature extraction and classification will achieve the same
 351 results. If $k > c$, (10) can obtain more than c projections to per-
 352 form feature selection or extraction, which provides the theoreti-
 353 cal guarantee for the performance of (10). From Propositions 1
 354 and 2, we can draw the following conclusion:

355 *Corollary 1:* If $\alpha \rightarrow 0$ and the matrix $\mathbf{X} \mathbf{X}^T$ is nonsingular,
 356 (1) and (3) have the same solution space.

358 C. Other Drawbacks of Ridge Regression

359 Adding L_2 -norm term for regression is of great importance to
 360 deal with the singular problem in (1). Moreover, it shows that no
 361 matter the matrix $\mathbf{X} \mathbf{X}^T$ is singular or nonsingular, the classical
 362 ridge regression in (3) is able to obtain the optimal solution and
 363 (1) is only a special case of (3) with the regularization parameter
 364 $\alpha = 0$. However, there are still some obvious disadvantages in

(3) since the optimal solution \mathbf{P}^* is not sparse, and thus it loses the feature selection function. Furthermore, the optimal solution \mathbf{P}^* in classical ridge regression model only contains the global information of the dataset and it ignores the local geometric structure. Thus, it is necessary to develop a new algorithm to deal with the above problems so as to enhance the effectiveness in feature extraction and pattern recognition. In the next section, we will propose a new model by jointing $L_{2,1}$ -norm regularization and locality regression to deal with the above problems.

IV. JOINTLY SPARSE LOCALITY REGRESSION ANALYSIS

In this section, the motivations and discussion are firstly present and then the proposed objective optimization problem as well as the local optimal solution will be presented.

A. The Motivations and Discussion

Based on the discussion in Section III-C and III-D, we can conclude the drawbacks of most existing regression methods into three aspects. First, due to the limitation of small-class problem, most regression methods cannot obtain enough projections to discover an effective projection matrix for discriminant feature extraction and classification. Second, the local structure of the data plays an important role in reconstructing the relationship between different data pairs in the low dimensional space. However, most existing regression methods do not take the local structure into consideration when performing feature selection or extraction. Third, there is no specific regression methods that are designed as regression form incorporating the local structures of the data as well as the sparsity of projections for feature selection and extraction.

Currently, deep learning technique is a research hotspot and it has been applied to the tasks of face recognition and object classification [65]. In spite of the high recognition rate of deep learning methods, behind is large-scale computing and long-term training. What is more, when the amount of data is not large enough, using deep learning methods for classification tends to obtain low performance because of the overfitting. In addition, most feature extraction methods based on deep learning [66], [67] do not consider the local structures of the data when doing convolutional operations. Even though they can obtain more abstract interpretation of the data, the relationship among different images is still missing. Therefore, developing efficient traditional feature extraction methods is still necessary for face recognition.

In conclusion, it is desirable to design a method that can solve the drawbacks of the existing regression methods and improve the performance of feature extraction to obtain high recognition rate with less computing time compared to the time-consuming and complicated deep learning methods.

B. The Objective Function of JSLR

To deal with the problems presented in Section III-C, Jointly Sparse Locality Regression Analysis (JSLR) is proposed to obtain a subset of jointly sparse projections for feature extraction and selection from the original data set. We also introduce the locality preserving regularized term to the model so as to characterize the local geometric structure of the data. Thus,

we present the objective function with joint $L_{2,1}$ -norm penalty and locality regularization. Let $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k]$ and $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k]$ be the variables of the following regression problem:

$$\begin{aligned} \bar{\mathbf{A}}, \bar{\mathbf{B}} = \arg \min_{\mathbf{A}, \mathbf{B}} & \left(\sum_{i=1}^n \|\mathbf{y}_i - \mathbf{x}_i^T \mathbf{B} \mathbf{A}^T\|_2^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \right. \\ & \left. + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|_2^2 \bar{\mathbf{w}}_{ij} \right) \\ & \mathbf{A}^T \mathbf{A} = \mathbf{I}, \end{aligned} \quad (14)$$

or in the matrix form

$$\begin{aligned} (\bar{\mathbf{A}}, \bar{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} & \left(\|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \right. \\ & \left. + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \bar{\mathbf{x}}_i - \mathbf{B}^T \bar{\mathbf{x}}_j\|_2^2 \bar{\mathbf{w}}_{ij} \right), \\ & \mathbf{A}^T \mathbf{A} = \mathbf{I}, \end{aligned} \quad (15)$$

where α_1 and α_2 are the regularization parameters. Since (14) and (15) have two variables and two kinds of norms in the model, they are not easy to be solved directly. Therefore, an alternatively iterative approach will be developed to solve the optimization problem in next section.

C. The Solutions of JSLR

From the definition of the $L_{2,1}$ -norm on the projection matrix \mathbf{B} , we have the diagonal matrix \mathbf{D}_B denoted as [54]

$$(\mathbf{D}_B)_{ii} = \frac{1}{2\|\mathbf{b}^i\|_2}, \quad (16)$$

where \mathbf{b}^i represents the i -th row of matrix \mathbf{B} .

Then from [56], we have the following equation:

$$\|\mathbf{B}\|_{2,1} = \text{Tr}(\mathbf{B}^T \mathbf{D}_B \mathbf{B}). \quad (17)$$

With the above preparation, we have

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 + \alpha_1 \|\mathbf{B}\|_{2,1} \\ & + \alpha_2 \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{B}^T \mathbf{x}_i - \mathbf{B}^T \mathbf{x}_j\|_2^2 \bar{\mathbf{w}}_{ij} \\ & = \text{Tr}(\mathbf{Y}^T \mathbf{Y} - 2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B} \\ & + \alpha_1 (\mathbf{B}^T \mathbf{D}_B \mathbf{B}) + \alpha_2 \mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B}). \end{aligned} \quad (18)$$

Since the optimization problem has two variables, we need to fix one to compute the other. For fixed \mathbf{A} , by setting the derivatives of (18) with respect to \mathbf{B} equaling to 0, (18) is minimized by

$$\bar{\mathbf{B}} = (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} \quad (19)$$

Hence, when \mathbf{A} is fixed, the objective function of 14 or 15 is minimized at the local optimal solution \mathbf{B} . When fixing \mathbf{B} , $\text{Tr}(\mathbf{Y}^T \mathbf{Y} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B} + \alpha_1 (\mathbf{B}^T \mathbf{D}_B \mathbf{B}) + \alpha_2 \mathbf{B}^T \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T \mathbf{B})$ becomes a constant and thus it can be

441 ignored. In such case, the following maximization problem gives
 442 the optimal solution to (18):

$$\max_{\mathbf{A}} \text{Tr}(\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A}) \text{ s.t. } \mathbf{A}^T \mathbf{A} = \mathbf{I}. \quad (20)$$

443 Let $\bar{\mathbf{A}}$ be the optimization of (20). From the Theorem 4 in [4],
 444 we have

$$\bar{\mathbf{A}} = \mathbf{U} \mathbf{V}^T, \quad (21)$$

445 where \mathbf{U}, \mathbf{V} is the SVD decomposition value of $\mathbf{Y}^T \mathbf{X}^T \mathbf{B}$.

446 In addition, we can also have the following conclusion from
 447 above formulation:

448 *Theorem 1:* Let $\bar{\mathbf{B}}$ be the local optimal solution of the optimi-
 449 zation problem (14) or (15). If $\alpha_1 \rightarrow 0$ and $\alpha_2 \rightarrow 0$, the linear
 450 subspace spanned by the optimal solution of (14) or (15) approxi-
 451 mates to the linear subspace spanned by \mathbf{P}^0 and \mathbf{P}^* , namely,
 452 $\text{span}(\bar{\mathbf{B}}) = \text{span}(\mathbf{P}^0)$ and $\text{span}(\bar{\mathbf{B}}) = \text{span}(\mathbf{P}^*)$.

453 *Proof:* The proof is in the Appendix.

454 For (19), when $\alpha_1 = 0$ and $\alpha_2 = 0$, then $\bar{\mathbf{B}} = (\mathbf{X} \mathbf{X}^T)^{-1}$
 455 $\mathbf{X} \mathbf{Y} \mathbf{A} = \mathbf{B}^* = \mathbf{P}^0 \mathbf{A}$, where \mathbf{P}^0 and \mathbf{B}^* is the optimal solu-
 456 tion corresponding to (1) and (10). As Proposition 1 and
 457 Proposition 2 have presented the relationship between (10)
 458 and (1), (10) and (3) respectively, it is easy for us to have the
 459 following conclusions: ■

460 *Corollary 2:* With the same assumptions and notations as in
 461 Theorem 1, when $\alpha_2 = 0, \alpha_1 \rightarrow 0$, \mathbf{P}^0 and $\bar{\mathbf{B}}$ have the same
 462 linear subspace, namely, $\text{span}(\mathbf{P}^0) = \text{span}(\bar{\mathbf{B}})$.

463 *Corollary 3:* With the same assumptions and notations as in
 464 Theorem 1, when $\alpha_1 = 0, \alpha_2 \rightarrow 0$, \mathbf{P}^0 and $\bar{\mathbf{B}}$ have the same
 465 linear subspace, namely, $\text{span}(\mathbf{P}^0) = \text{span}(\bar{\mathbf{B}})$.

466 In summary, from Theorem 1, Corollary 2 and Corollary 3,
 467 we can know that either (14) or (15) provides a basic theoretical
 468 guarantee for the effectiveness of the proposed regression
 469 model. Namely, when the parameters of the proposed model are
 470 set suitably, the optimal solution space of the ridge regression
 471 can be derived from (15). This means that the optimal projection
 472 of JSLR can approximate to the subspace spanned by the tradi-
 473 tional regression models. Besides, by utilizing the advantages
 474 of $L_{2,1}$ -norm regularization and locality preserving property, the
 475 proposed model is able to compute the jointly sparse projections
 476 and preserve the local geometric structure of the data for feature
 477 extraction. The detail of the iterative algorithm was illustrated
 478 in Algorithm 1.

479 D. Comparison and Discussion

480 In this section, we compare our algorithm JSLR with other
 481 methods, such as PCA, SPCA, LDA, LPP and so on. Both PCA
 482 and SPCA are outstanding in data processing and dimensionality
 483 reduction. PCA projects the original d -dimensional data onto
 484 $k (< d)$ -dimensional linear subspace with the combination of
 485 all the original variables. SPCA aims to produce modified sparse
 486 principal components by lasso (or elastic net) technique. But it
 487 just focuses on the global structure of the original data and ignore
 488 the local structure. Different from SPCA, JSLR can efficiently
 489 preserve the local geometric structure of the data set.

490 Some other subspace learning algorithms, LPP, NPE, etc. are
 491 able to preserve local structure of the original data. However,
 492 they cannot provide the jointly sparse property for the learned

Algorithm 1: JSLR Algorithm

Input: The training data $\mathbf{X} \in R^{d \times n}$,

the training data label $\mathbf{Y} \in R^{n \times c}$,

matrices $\bar{\mathbf{D}} \in R^{n \times n}$, $\bar{\mathbf{W}} \in R^{n \times n}$,

the objective dimension k ($k = 1, 2, \dots, n$),

maximum number of the iteration: maxStep .

Step 1: Compute matrices $\bar{\mathbf{D}}, \bar{\mathbf{W}}$, and initialize matrix \mathbf{D}_B ,
 $step = 0$, $\text{converged} = \text{false}$.

Step 2: While !converged and $step <= \text{maxStep}$

- Compute \mathbf{B} using

$$\mathbf{B} = (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}$$

- Compute \mathbf{D}_B using $(\mathbf{D}_B)_{ii} = \frac{1}{2\|\mathbf{b}^i\|_2}$

- Compute \mathbf{A} using $\mathbf{A} = \mathbf{U} \mathbf{V}^T$

- Update $\text{converged} = \text{true}$ when \mathbf{B} is approximately
 changeless.

Step 3: Standardize the matrix \mathbf{B} to a final normalized
 matrix

and return it for feature selection.

Output: Low-dimensional discriminative subspace

$$\mathbf{B} \in R^{d \times k}, k = 1, 2, \dots, n.$$

493 subspace. Compared with them, JSLR achieves this goal by
 494 adding $L_{2,1}$ -norm regularization to make the elements in some
 495 rows of the projection to be 0 for efficient feature extraction and
 496 selection.

497 Ridge regression is frequently used in face recognition. How-
 498 ever, when the class number of training sample is too small, ridge
 499 regression cannot obtain more projections than the number of
 500 the classes for feature extraction. The same problem exists in
 501 LDA. In contrast, the number of the projections of JSLR is not
 502 limited by the number of the classes in training data. In spite
 503 of given a small number of classes in training sample set, JSLR
 504 can obtain any number of projections for feature selection and
 505 the number of the projections is freely set by the users.

506 In summary, the advantages of JSLR against PCA, SPCA,
 507 LDA, ridge regression and LPP are that JSLR can obtain joint
 508 sparsity and preserve the local structure for pattern recogni-
 509 tion. Another major difference between JSLR and other clas-
 510 sical methods is that the number of the training sample classes
 511 in JSLR is allowed to be very small but it still can learn more
 512 projections than the number of classes. These advantages make
 513 JSLR achieve high recognition rate.

V. THEORETICAL ANALYSIS

514 In this section, we present the theoretical analysis including
 515 convergence analysis and computational complexity analysis.

A. The Convergence

516 To verify the convergences of the proposed iterative algo-
 517 rithm, we begin with the following Lemmas:

518 *Lemma 1:* [54] For any two non-zero constants a and b , we
 519 have the following inequality:

$$\sqrt{a} - \frac{a}{2\sqrt{b}} \leq \sqrt{b} - \frac{b}{2\sqrt{b}}. \quad (22)$$

522 Lemma 2: [54] Denoted \mathbf{V} as any nonzero matrix, $\mathbf{V} \in R$,
 523 the following inequality holds:

$$\sum_i \|\mathbf{v}_t^i\|_2 - \sum_i \frac{\|\mathbf{v}_t^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2} \leq \sum_i \|\mathbf{v}_{t-1}^i\|_2 - \sum_i \frac{\|\mathbf{v}_{t-1}^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2}, \quad (23)$$

524 where $\mathbf{v}_t^i, \mathbf{v}_{t-1}^i$ denote the i -th row of matrix \mathbf{V}_t and \mathbf{V}_{t-1} .

525 Proof: Let $\|\mathbf{v}_t^i\|_2^2$ and $\|\mathbf{v}_{t-1}^i\|_2^2$ be the substitute of a and b
 526 in (22), the following inequality is valid for any i .

$$\|\mathbf{v}_t^i\|_2 - \frac{\|\mathbf{v}_t^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2} \leq \|\mathbf{v}_{t-1}^i\|_2 - \frac{\|\mathbf{v}_{t-1}^i\|_2^2}{2\|\mathbf{v}_{t-1}^i\|_2}, \quad (24)$$

527 Thus, (23) as the sum form of (24) also holds (22). With
 528 the above Lemma 1 and Lemma 2, we have the following
 529 theorem: ■

530 Theorem 2: Given all the parameters in the objective function
 531 except \mathbf{A} and \mathbf{B} , the iterative approach shown in Algorithm
 532 1 will monotonically decrease the objective function value of
 533 (14) or (15) in each iteration and provides a local optimal solution
 534 of the problem.

535 Proof: For simplicity, we denote the objective function of
 536 (18) as $F(\mathbf{B}, \mathbf{A}) = F(\mathbf{B}, \mathbf{A}, \mathbf{D}_B)$. Suppose for the $(t-1)$ -th
 537 iteration, both \mathbf{A}_{t-1} and \mathbf{B}_{t-1} can be obtained. Then we have
 538 the following inequality from (19):

$$F(\mathbf{B}_t, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}) \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (25)$$

539 For \mathbf{A}_t , as its optimal value comes from SVD and this will
 540 further decrease the value of the objective function, it goes

$$F(\mathbf{B}_t, \mathbf{A}_t, (\mathbf{D}_B)_{t-1}) \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (26)$$

541 In (18), since $\mathbf{Y}^T \mathbf{Y}$ is a constant, it can be ignored and we need
 542 to minimize

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}) \end{aligned}$$

543 As we have obtained the optimal \mathbf{B}_t and \mathbf{A}_t , then the following
 544 inequality holds:

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) \quad (27) \end{aligned}$$

545 That is

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) + \alpha_1 \sum_i \frac{\|\mathbf{b}_t^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) + \alpha_1 \sum_i \frac{\|\mathbf{b}_{t-1}^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \quad (28) \end{aligned}$$

TABLE I
THE COMPUTATIONAL COMPLEXITIES

Iteration variable	computational complexities
\mathbf{B}	$O(d^3)$
\mathbf{D}_B	$O(d^2)$
\mathbf{A}	$O(d^3)$

546 Then the above inequality indicates

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) \\ & + \alpha_1 \sum_i \|\mathbf{b}_t^i\|_2 - \alpha_1 \left(\sum_i \|\mathbf{b}_t^i\|_2 - \sum_i \frac{\|\mathbf{b}_t^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \right) \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) \\ & + \alpha_1 \sum_i \|\mathbf{b}_{t-1}^i\|_2 - \alpha_1 \left(\sum_i \|\mathbf{b}_{t-1}^i\|_2 - \sum_i \frac{\|\mathbf{b}_{t-1}^i\|_2^2}{2\|\mathbf{b}_{t-1}^i\|_2} \right) \quad (29) \end{aligned}$$

547 According to Lemma 2, we further have

$$\begin{aligned} & \text{Tr}(-2\mathbf{B}_t^T \mathbf{X} \mathbf{Y} \mathbf{A}_t + \mathbf{B}_t^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_t) + \alpha_1 \sum_i \|\mathbf{b}_t^i\|_2 \\ & \leq \text{Tr}(-2\mathbf{B}_{t-1}^T \mathbf{X} \mathbf{Y} \mathbf{A}_{t-1} + \mathbf{B}_{t-1}^T (\mathbf{X} \mathbf{X}^T + \alpha_1 (\mathbf{D}_B)_{t-1} \\ & + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T) \mathbf{B}_{t-1}) + \alpha_1 \sum_i \|\mathbf{b}_{t-1}^i\|_2. \quad (30) \end{aligned}$$

548 That is

$$\begin{aligned} & F(\mathbf{B}_t, \mathbf{A}_t) = F(\mathbf{B}_t, \mathbf{A}_t, (\mathbf{D}_B)_t) \\ & \leq F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}) = F(\mathbf{B}_{t-1}, \mathbf{A}_{t-1}, (\mathbf{D}_B)_{t-1}). \quad (31) \end{aligned}$$

549 From (31), we can conclude that the objective function value
 550 of (14) or (15) is monotonically decreased via the updating rule
 551 presented in Algorithm 1. Therefore, the proposed iterative al-
 552 gorithm finally converges to the local optimal solution. ■

B. Computational Complexity Analysis

553 For simplicity, we assume the dimension of training samples
 554 is d . Our proposed algorithm aims to compute the matrix \mathbf{A} and
 555 \mathbf{B} . Computing \mathbf{B} in (19) needs $O(d^3)$ while computing \mathbf{D}_B
 556 in (16) needs $O(d^2)$. Since SVD of $\mathbf{Y}^T \mathbf{X}^T \mathbf{B}$ also needs $O(d^3)$,
 557 then the computational complexity of \mathbf{A} is also $O(d^3)$. It is easy
 558 to know that the main complexity of the algorithm is $O(Td^3)$,
 559 where T denotes the number of iterations for convergence. Ta-
 560 ble I lists the computational complexities of each variable. 561

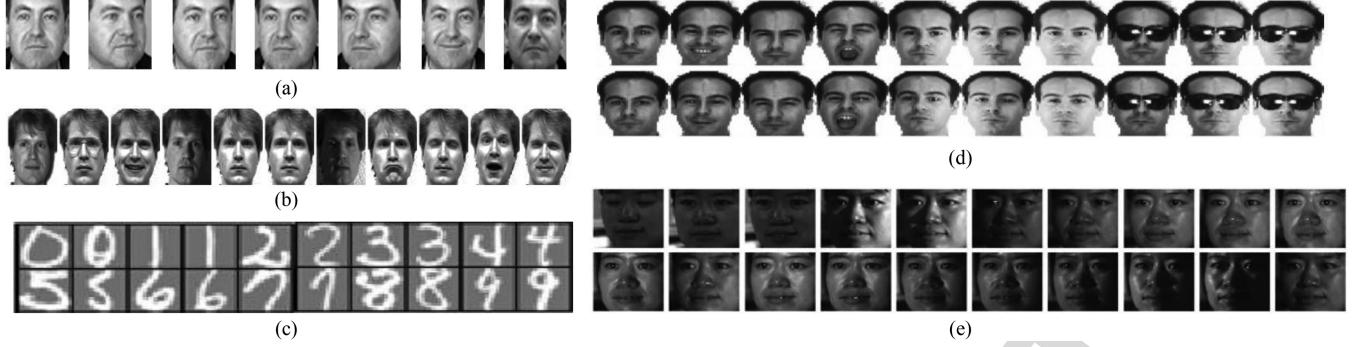


Fig. 1. Examples from FERET, Yale, USPS, AR, and CMU PIE data sets. (a) FERET. (b) Yale. (c) USPS. (d) AR. (e) CMU PIE.

VI. EXPERIMENTS

In this section, to evaluate the proposed JSLR algorithm for feature selection, we conducted a set of experiments from three aspects: experiments on small-scale face databases, experiments on large-scale databases and experiments based on deep learning. In experiments, several classical as well as state-of-the-art methods are used as compared methods. They are the classical principle component analysis method PCA, the classical Ridge Regression (RR) [60], the Linear Discriminant Analysis (LDA) [5], the traditional sparse learning method SLDA based on L_1 -norm [23], the local structure learning method Locality Preserving Projection (LPP) [32], the regression analysis of locality preserving projections via sparse penalty (SpLPP) [61] which applies sparsity penalty and minimization based on L_1 -norm to locality preserving projections, the dictionary learning methods (i.e. label consistent K-SVD (LC-KSVD2) [68] and the Locality Constrained and Label Embedding Dictionary Learning (LCLE-DL) [69]), the most related $L_{2,1}$ -norm regularization methods for feature selection and subspace learning (i.e. Unsupervised Discriminative Feature Selection (UDFS) [56] and Robust Feature Selection (RFS) [54]). In addition, the proposed method without $L_{2,1}$ -norm regularization named JSLR($\alpha_1 = 0$) (i.e. the second term in the proposed objective function in Eq. (14) is removed) was added as a compared method to all experiments to evaluate the effectiveness of the jointly sparse regularization.

In all experiments we make comparison in the avenue of deep learning (the method is called Deep-NN in this paper). Deep-NN is completed by the following two steps. Firstly, we use the deep convolutional neural network (CNN) as the feature extractor to obtain the deep features of all samples. This process is similar to [70]. Secondly, we use the nearest neighbor classifier (NN) for classification. For the proposed JSLR, we also use the deep features instead of the traditional image features as input and this method is called Deep-JSLR for easy understanding. Note that the deep features of character database are obtained according to the tutorial of MNIST network on official Caffe site (<http://caffe.berkeleyvision.org/gathered/examples/mnist.html>.)

A. Experiments on Small-Scale Database

In this section, experiments on four databases, including FERET, AR, CMU PIE and Yale database, were conducted to

evaluate the performance of the proposed method versus the compared methods under different variations of facial expression and lighting condition.

1) *Experiments on FERET Face Database*: The FERET face database [71] includes 1,400 images of 200 individuals (each individual has seven images). In the experiment, the facial portion of each original image was automatically cropped based on the location of the eyes, and the cropped images were resized to 40×40 pixels. The sample images of one person are shown in Fig. 1(a).

Experimental Setting: For all the databases, the image set is partitioned into two parts, i.e. the gallery and probe sets. In each database, l (l is no more than the number of class) images of each class are randomly selected to form the gallery set and the remaining images are used as the probe set. PCA was used as pre-processing to reduce the dimension of data. Then the proposed method and the compared methods were used to perform feature extraction, independently. Finally, nearest neighbor classifier was used for classification. The experiments were independently performed 10 times. The average recognition rates and the corresponding dimensions as well as the standard deviations of each method were listed on the Table II. Besides, the comparison results were also shown in the Fig. 2(c)–(f) when 5 images of each individual were randomly selected for training and the remaining images were used for testing. The dimensions of the projection matrices were set as empirical value and marked on the horizontal axis. The variables except parameter α_1 and α_2 in JSLR were randomly initialized in our experiments.

Exploration of the Performance of the Parameters: In order to explore the optimal parameters for JSLR on different data sets, we analyzed the values of the parameters α_1 (Alpha1) and α_2 (Alpha2). For the other compared methods, since in most of cases the best performance lie on the area of $[10^{-3}, 10^3]$, as introduced in the corresponding papers, we fixed their parameters on the area of $[10^{-3}, 10^3]$ and report the best results.

In this experiment, we analyze the impacts of various parameter values on the performance of JSLR and the average recognition rates of different dimensions from 5 to 200. Table II shows the best average recognition rates based on 10 times running and the corresponding dimensions as well as the standard deviations of each method with l ($l = 4, 5$) images of each individual for training while the remaining images were used for testing.

TABLE II
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF DIFFERENT METHODS ON FERET FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	54.55	54.72	60.77	46.80	61.23	53.12	63.82	55.53	59.42	46.72	73.22	74.15	99.47	100.00
	± 8.54	± 8.74	± 20.04	± 9.61	± 19.96	± 12.55	± 21.07	± 23.33	± 13.46	± 9.55	± 18.77	± 18.34	± 0.46	± 0.00
	28^*5	27^*5	30^*5	30^*5	30^*5	12^*5	40^*5	13^*5	30^*5	30^*5	30^*5	29^*5	30^*5	26^*5
5	65.50	65.90	76.45	61.50	77.18	66.65	78.35	74.92	67.95	60.08	90.70	91.35	99.55	100.00
	± 5.48	± 5.65	± 7.76	± 4.21	± 7.61	± 5.27	± 7.33	± 10.55	± 6.33	± 6.69	± 6.67	± 6.42	± 0.33	± 0.00
	30^*5	27^*5	30^*5	30^*5	29^*5	11^*5	40^*5	8^*5	29^*5	30^*5	27^*5	29^*5	30^*5	20^*5

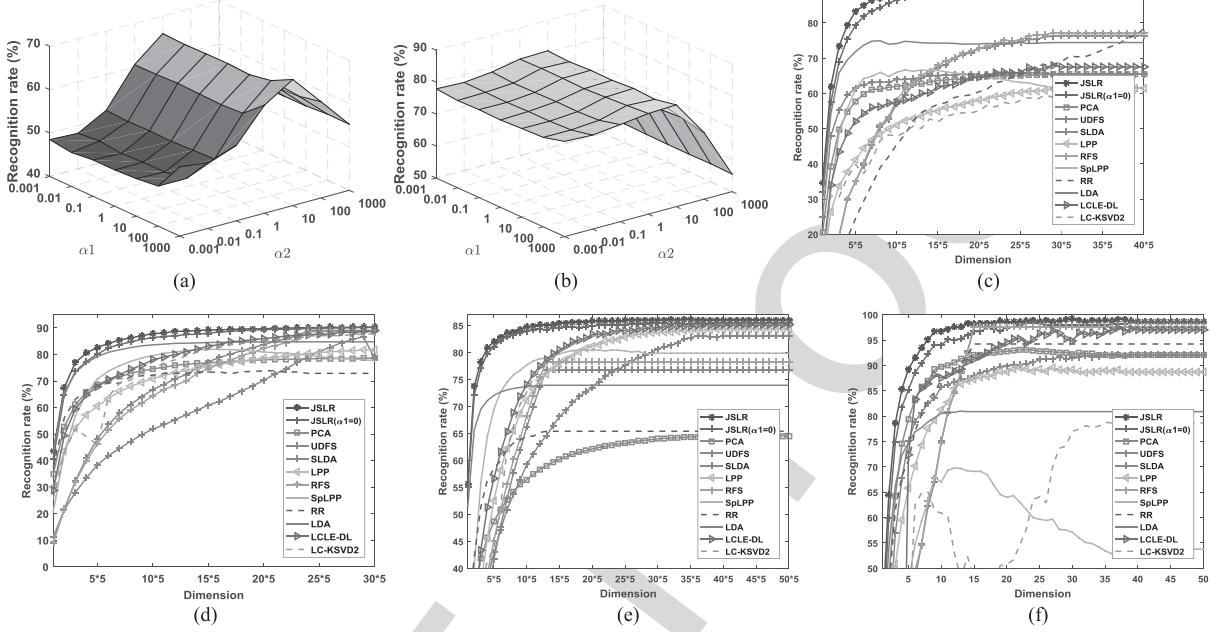


Fig. 2. The recognition rate versus the parameters α_1 and α_2 on the (a) FERET and (b) AR face database, respectively. The recognition rates (%) versus the dimensions of different methods on the (c) FERET, (d) AR, (e) CMU PIE, (f) Yale face databases, respectively.

Fig. 2(a) shows the recognition rates when the two parameters α_1 and α_2 change from 10^{-3} to 10^3 . Fig. 2(c) shows the average recognition rates versus various dimensions of different methods.

It is easy to know that the optimal value of the parameter α_1 lies on the area of $[10^{-3}, 10^2]$ while the optimal value of the parameter α_2 lies on the area of $[10^{-2}, 10^3]$. In other words, JSLR is efficient and robust among these areas. By contrast, when the values of the two parameters lie on other areas, it will cause the larger decline of the recognition rates.

As it can be seen from Fig. 2(c), the recognition rates of JSLR as well as $\text{JSLR}(\alpha_1=0)$ are the highest. The results shown in Table II and Fig. 2(c) indicate that JSLR and $\text{JSLR}(\alpha_1=0)$ outperform PCA, UDFS, SLDA, RFS, LPP and SpLPP, RR, LDA, LCLE-DL, LC-KSVD2 in feature extraction. Besides, from Table II, we can easily know that Deep-JSLR outperforms Deep-NN.

2) *Experiments on AR Face Database*: The AR face database [72] contains the pictures of 120 individuals (each individual has 20 images). The face portion of each image was manually cropped (because of missing eye coordinates) and then normalized to 50×40 pixels. The sample images of one person are shown in Fig. 1(d).

In this experiment, we randomly selected l ($l = 4, 5, 6$) images of each individual for training, and the rest of the images in the data set were used for testing. From Fig. 2(b), we can know that the optimal values of parameter α_1 and α_2 were both $[10^{-3}, 10^2]$. Thus, we used this area for JSLR to obtain the comparison results. Table III listed the performance of different methods. Fig. 2(d) showed the average testing recognition rates. It is obvious that JSLR or Deep-JSLR outperforms the other methods.

3) *Experiments on CMU PIE Database*: The CMU PIE face database [73] contains 68 individuals with 41,368 face images as a whole. We selected a subset (C29) containing 1632 images from 68 individuals (each providing 24 images). All of these face images were automatically aligned based one-eye coordinates and cropped to 32×32 pixels. Fig. 1(e) shows the sample images from this database.

In this experiment, l ($l = 4, 5, 6$) images of each individual were randomly selected for training, and the rest of the images in the data set were used for testing. The optimal areas of α_1 and α_2 were the same with the areas on AR database. Table IV presents the performance of different methods. Fig. 2(e) shows the average testing recognition rates and indicates that JSLR outperforms the other methods again.

TABLE III
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON AR FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	76.86	85.66	87.97	80.40	87.41	79.03	68.29	80.06	89.89	84.97	88.82	89.08	83.43	90.23
	± 4.56	± 8.48	± 10.95	± 9.25	± 11.41	± 7.76	± 5.89	± 11.78	± 7.76	± 6.96	± 10.79	± 10.64	± 12.78	± 9.77
	$30*5$	$29*5$	$24*5$	$30*5$	$24*5$	$12*5$	$21*5$	$23*5$	$29*5$	$30*5$	$30*5$	$28*5$	$29*5$	$28*5$
5	78.67	86.73	89.07	81.91	88.58	82.01	73.91	84.78	88.87	87.88	90.19	90.33	89.26	94.70
	± 5.41	± 8.89	± 11.27	± 9.78	± 11.67	± 8.34	± 7.99	± 11.77	± 8.28	± 7.44	± 10.63	± 10.60	± 9.36	± 7.02
	$30*5$	$29*5$	$24*5$	$30*5$	$24*5$	$15*5$	$21*5$	$21*5$	$30*5$	$30*5$	$29*5$	$29*5$	$29*5$	$29*5$
6	80.47	90.35	94.24	86.11	93.85	86.76	79.41	91.82	91.50	92.26	95.05	95.23	93.27	97.52
	± 4.97	± 5.85	± 8.11	± 7.89	± 8.51	± 6.15	± 6.02	± 9.38	± 5.92	± 5.85	± 8.02	± 8.04	± 3.07	± 1.84
	$30*5$	$28*5$	$24*5$	$30*5$	$24*5$	$16*5$	$21*5$	$22*5$	$30*5$	$30*5$	$28*5$	$28*5$	$29*5$	$30*5$
15	92.08	92.08	99.22	98.80	99.22	95.03	96.30	99.12	98.02	97.10	99.55	99.60	97.68	99.73
	± 18.01	± 3.44	± 20.38	± 0.71	± 0.41	± 2.54	± 1.57	± 0.68	± 1.16	± 1.26	± 0.32	± 0.32	± 2.91	± 0.74
	$30*5$	$28*5$	$24*5$	$30*5$	$24*5$	$16*5$	$21*5$	$21*5$	$30*5$	$30*5$	$27*5$	$19*5$	$30*5$	$30*5$
16	92.35	92.35	99.27	98.92	99.27	95.19	96.69	99.25	98.08	97.50	99.63	99.77	98.67	99.98
	± 18.07	± 3.48	± 20.41	± 0.78	± 0.45	± 1.72	± 1.90	± 0.53	± 1.17	± 1.01	± 0.28	± 0.25	± 1.38	± 0.07
	$30*5$	$28*5$	$24*5$	$30*5$	$24*5$	$16*5$	$21*5$	$21*5$	$30*5$	$30*5$	$18*5$	$19*5$	$30*5$	$22*5$

TABLE IV
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON CMU PIE FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	57.11	77.62	69.25	80.01	71.55	76.21	56.33	36.31	81.43	81.05	80.60	81.43	96.32	97.54
	± 12.66	± 8.85	± 13.46	± 9.16	± 12.92	± 8.32	± 10.50	± 8.54	± 6.78	± 8.15	± 11.12	± 10.4	± 1.10	± 0.92
	$40*5$	$37*5$	$13*5$	$40*5$	$13*5$	$20*5$	$13*5$	$13*5$	$39*5$	$36*5$	$34*5$	$30*5$	$38*5$	$37*5$
5	64.52	83.13	76.76	84.35	78.25	80.48	65.39	73.92	85.11	85.22	85.44	86.08	96.73	97.96
	± 12.27	± 4.65	± 11.02	± 4.85	± 10.19	± 6.86	± 9.26	± 10.23	± 4.32	± 3.64	± 5.19	± 4.54	± 1.25	± 1.17
	$40*5$	$36*5$	$13*5$	$40*5$	$13*5$	$19*5$	$13*5$	$13*5$	$37*5$	$39*5$	$33*5$	$34*5$	$39*5$	$34*5$
6	68.84	85.63	78.45	85.62	79.77	83.51	71.36	78.19	85.50	85.28	85.16	85.85	96.99	98.14
	± 10.20	± 4.74	± 9.59	± 5.04	± 7.85	± 5.44	± 10.07	± 10.30	± 3.56	± 3.80	± 5.15	± 4.54	± 1.10	± 1.07
	$40*5$	$38*5$	$13*5$	$37*5$	$13*5$	$19*5$	$13*5$	$13*5$	$40*5$	$35*5$	$38*5$	$33*5$	$39*5$	$38*5$
19	94.24	94.24	92.29	92.09	92.09	91.88	94.41	91.06	87.00	90.09	92.91	93.12	99.32	100.00
	± 5.91	± 5.91	± 8.22	± 8.03	± 8.34	± 8.39	± 5.46	± 9.55	± 9.74	± 7.52	± 7.65	± 7.55	± 0.77	± 0.00
	$40*5$	$38*5$	$13*5$	$40*5$	$13*5$	$19*5$	$13*5$	$13*5$	$40*5$	$40*5$	$40*5$	$33*5$	$33*5$	$14*5$
20	93.16	93.16	90.29	90.18	90.11	90.22	93.68	89.04	85.57	88.54	91.29	91.47	99.34	100.00
	± 7.17	± 7.17	± 10.25	± 9.95	± 10.28	± 9.95	± 6.56	± 11.73	± 11.96	± 9.17	± 9.44	± 9.28	± 0.76	± 0.00
	$40*5$	$38*5$	$13*5$	$40*5$	$13*5$	$19*5$	$13*5$	$13*5$	$40*5$	$40*5$	$23*5$	$30*5$	$33*5$	$14*5$

TABLE V
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON YALE FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha_1=0$)	JSLR	Deep-NN	Deep-JSLR
4	91.24	90.48	96.95	85.05	96.48	68.10	91.62	66.76	97.90	78.86	97.71	98.95	99.71	100.00
	± 3.46	± 4.02	± 2.49	± 5.06	± 2.58	± 9.27	± 3.24	± 10.23	± 1.41	± 41.68	± 1.61	± 1.05	± 0.40	± 0.00
	21	39	15	37	15	11	15	10	40	33	31	27	40	12
5	93.11	92.11	97.56	89.56	97.56	69.78	94.22	81.00	97.44	78.78	98.44	99.22	99.56	100.00
	± 3.60	± 4.57	± 2.39	± 3.58	± 2.21	± 5.49	± 3.57	± 5.50	± 1.80	± 41.78	± 1.57	± 1.20	± 0.47	± 0.00
	22	36	15	22	15	12	15	13	37	35	39	30	40	11

692 4) *Experiments on Yale Database*: The Yale face database
693 [43] contains 165 grayscale images of 15 individuals. Each
694 image was manually cropped (because of no eye coordinates
695 provided) and resized to 50×40 pixels. Fig. 1(b) shows the
696 sample images from this database.
697 In this experiment, l ($l = 4, 5$) images of each individual were
698 randomly selected for training, and the rest of the images in the
699 data set were used for testing. The values of α_1 and α_2 were both
700 from 10^{-3} to 10^2 . The performances of the different methods
701 are shown in Table V. Fig. 2(f) shows the average recognition
702 rates. It clearly indicates that JSLR and Deep-JSLR can obtain
703 the best performance when the traditional image features and
704 deep features are used as input.

705 B. Experiments on Large-Scale Database

706 In this section, two different databases are used to evaluate
707 the performance of the proposed method based on large-scale
708 data learning.

709 1) *Experiments on USPS Database*: The United States
710 Postal Service(USPS) database [55] consists of 1,100 of each
711 handwritten digit (0-9). The images in this database are resized
712 to 16×16 pixels. The performance of JSLR on large-sample
713 data set was evaluated on this database. Fig. 1(c) shows the sample
714 images from this database.

715 In this experiment, l ($l = 400, 500, 600, 5, 10$) images of each
716 class were randomly selected for training, and the rest of the images
717 in the data set were used for testing. The parameter α_1 and
718 α_2 were both from 10^{-3} to 10^3 . Table VI shows the performance
719 of the different methods. From the result, we can know that JSLR
720 or Deep-JSLR can achieve better performance than other compared
721 methods on this database. Particularly, when only 5/1100 images
722 of each class are used for training, the proposed method can
723 obtain higher recognition rate than not only the compared methods
724 but also the Deep-NN.

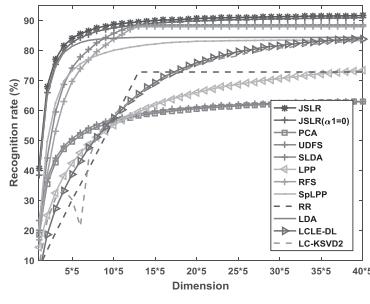
725 2) *Experiments on PIE67 \times 170 Database*: The PIE67 \times
726 170 database is a subset of the CMU PIE face database [73].

TABLE VI
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON USPS FACE DATABASE

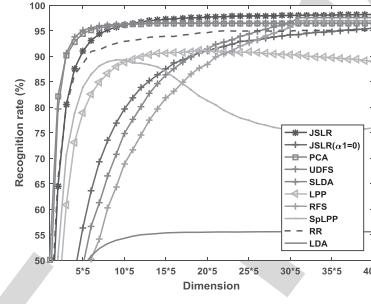
Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR	Deep-NN	Deep-JSLR
400	90.55	91.17	91.37	89.41	90.42	91.02	90.34	86.78	87.21	90.78	92.02	92.35	91.44	95.74
	± 0.13	± 0.12	± 0.35	± 0.30	± 0.40	± 0.19	± 0.17	± 0.00	± 0.56	± 0.19	± 0.22	± 0.18	± 0.28	± 0.27
	7*5	10*5	2*5	6*5	2*5	6*5	2*5	1*5	4*5	30*5	30*5	20*5	1*5	15*5
500	91.02	91.75	91.91	90.14	90.81	91.90	90.87	87.23	87.22	91.60	92.52	92.77	92.18	96.00
	± 0.15	± 0.18	± 0.23	± 0.21	± 0.31	± 0.16	± 0.23	± 0.00	± 0.59	± 0.27	± 0.18	± 0.14	± 0.17	± 0.26
	7*5	9*5	2*5	7*5	2*5	6*5	2*5	1*5	4*5	30*5	30*5	28*5	1*5	21*5
600	91.49	92.20	92.28	90.58	91.40	92.51	91.45	87.49	87.15	92.15	92.82	93.08	92.67	96.09
	± 0.17	± 0.24	± 0.35	± 0.17	± 0.45	± 0.13	± 0.25	± 0.00	± 0.71	± 0.33	± 0.28	± 0.31	± 0.26	± 0.21
	7*5	8*5	2*5	6*5	2*5	6*5	2*5	1*5	4*5	30*5	28*5	25*5	1*5	20*5
5	61.43	61.43	56.83	49.41	56.23	15.76	57.19	61.45	62.00	62.04	59.51	67.83	56.96	66.96
	± 2.08	± 2.08	± 2.30	± 3.10	± 2.43	± 2.28	± 2.24	± 3.53	± 2.40	± 2.34	± 2.24	± 2.39	± 2.28	± 1.72
	7*5	8*5	2*5	6*5	2*5	6*5	2*5	1*5	4*5	30*5	8*5	10*5	$100*5$	10*5
10	69.86	69.86	55.78	49.49	54.73	14.14	55.64	61.02	68.87	67.81	58.67	76.22	65.66	78.15
	± 2.18	± 2.18	± 2.15	± 2.64	± 2.00	± 2.30	± 2.05	± 2.16	± 1.75	± 1.91	± 1.87	± 2.64	± 1.68	± 1.64
	7*5	8*5	2*5	6*5	2*5	6*5	2*5	1*5	4*5	30*5	12*5	20*5	$100*5$	20*5



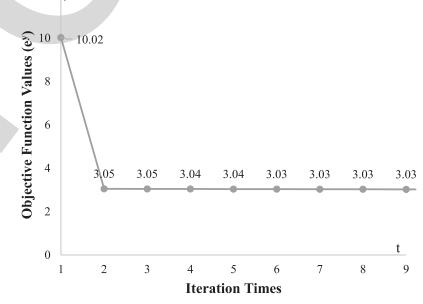
(a)



(b)



(c)



(d)

Fig. 3. (a) Sample images on the LFW database. The recognition rates (%) versus the dimensions of different methods on the (b) PIE67 \times 170, (c) LFW databases, respectively. (d) An example of the convergence curve of JSLR on Yale database.

TABLE VII
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON PIE67 \times 170 FACE DATABASE

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR	Deep-NN	Deep-JSLR
10	45.75	45.75	78.63	59.72	76.14	69.71	58.85	72.74	82.72	78.78	81.47	83.64	90.90	95.62
	± 3.71	± 3.72	± 3.55	± 3.33	± 4.03	± 3.69	± 3.01	± 2.70	± 2.86	± 2.11	± 2.93	± 3.17	± 1.21	± 1.31
	40*5	40*5	13*5	40*5	13*5	21*5	13*5	13*5	40*5	40*5	40*5	40*5	40*5	40*5
20	62.89	62.89	88.41	73.32	87.84	83.38	72.75	84.29	86.84	86.19	90.84	91.55	93.98	97.70
	± 3.05	± 3.06	± 2.09	± 3.12	± 2.11	± 3.09	± 2.66	± 2.11	± 1.42	± 1.23	± 1.42	± 1.37	± 1.07	± 0.73
	40*5	40*5	13*5	40*5	13*5	30*5	13*5	13*5	40*5	40*5	40*5	40*5	40*5	40*5
30	69.85	69.85	91.46	77.99	91.05	87.42	83.74	87.98	87.81	87.96	93.38	93.70	95.36	98.38
	± 3.57	± 3.58	± 2.22	± 2.76	± 2.30	± 2.72	± 3.13	± 2.78	± 1.61	± 1.27	± 1.82	± 1.78	± 1.00	± 0.67
	40*5	40*5	13*5	40*5	13*5	36*5	13*5	13*5	40*5	39*5	40*5	40*5	40*5	40*5

727 There are total 11,390 images from 67 individuals and each
728 individual has 170 images on this database. The experiment on
729 this database is conducted to evaluate the performance of JSLR
730 as well as the compared methods on the occasion when there are
731 various facial expression, lighting condition and angle on the
732 face images.

733 In this experiments, l ($l = 10, 20, 30$) images of each individual are randomly selected for training and the remaining are used for testing. The recognition rates of all methods are shown

734 in Fig. 3(b) and Table VII. From Fig. 3(b), JSLR as well as
735 JSLR($\alpha=0$) obtain higher recognition rate than other methods, which indicates that the proposed method is superior to
736 other methods even without the $L_{2,1}$ -norm regularization term
737 (this can also be verified by Fig. 2(c) and Table II).
738

C. Experiments Based on Deep Learning

739 In this section, experiments on three database (AR, the standard
740 subsets of the FERET and the LFW databases [74]) were
741

TABLE VIII
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON AR FACE DATABASE BASED ON DEEP LEARNING

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	Deep-NN	JSLR ($\alpha=0$)	JSLR
4	87.12	87.12	87.53	83.46	88.80	87.12	74.53	82.37	90.01	89.89	83.43	13.49	90.23
	± 12.26	± 12.33	± 9.26	± 10.96	± 10.73	± 11.90	± 9.91	± 8.81	± 6.50	± 6.89	± 12.78	± 7.47	± 9.77
	30*5	30*5	24*5	22*5	24*5	30*5	24*5	23*5	30*5	30*5	29*5	30*5	28*5
5	91.84	91.84	91.16	89.82	93.17	91.84	80.74	89.63	94.67	94.10	89.26	76.91	94.70
	± 9.36	± 9.43	± 7.84	± 8.20	± 7.79	± 8.83	± 7.82	± 7.56	± 4.61	± 4.60	± 9.36	± 5.51	± 7.02
	30*5	30*5	24*5	20*5	24*5	30*5	24*5	21*5	30*5	30*5	29*5	30*5	29*5
6	95.40	95.33	94.39	93.99	96.16	95.33	84.64	93.85	96.81	96.67	93.27	92.93	97.52
	± 2.77	± 2.75	± 3.84	± 3.14	± 2.74	± 2.82	± 6.58	± 3.65	± 2.25	± 1.78	± 3.07	± 3.58	± 1.84
	29*5	30*5	24*5	16*5	24*5	30*5	24*5	18*5	30*5	30*5	29*5	30*5	30*5

TABLE IX
THE MAXIMAL RECOGNITION RATE OF ALL METHODS ON THE Fb, Fc, Dup1, Dup2 FACE DATABASE BASED ON DEEP LEARNING

Algorithm	Fb (dim=216, 512)	Fc (dim=216, 512)	Dup1(dim=216, 512)	Dup2(dim=216, 512)
PCA	99.41, 99.41	99.48, 99.48	98.20, 98.20	98.29, 98.29
UDFS	80.75, 99.41	65.46, 99.48	43.07, 98.20	45.30, 98.29
SLDA	98.49, 98.91	100.00, 100.00	84.63, 88.50	89.74, 92.31
LPP	98.58, 99.16	90.03, 92.66	86.75, 91.45	86.75, 91.45
RFS	99.58, 99.58	100.00, 100.00	97.51, 98.48	98.72, 99.15
SpLPP	99.41, 99.50	98.97, 99.48	97.78, 98.06	98.72, 99.15
RR	98.74, 98.91	98.45, 98.45	94.46, 94.46	96.58, 96.58
LDA	99.41, 99.41	99.48, 99.48	98.20, 98.20	98.29, 98.29
LCLE-DL*	-	-	-	-
LC-KSVD2*	-	-	-	-
Deep-NN	99.41, 99.50	99.48, 99.48	98.48, 98.48	98.72, 98.72
JSLR ($\alpha=0$)	98.58, 99.41	99.48, 100.00	86.01, 93.77	88.03, 94.44
JSLR	99.67, 99.67	100.00, 100.00	98.75, 98.89	99.57, 99.57

*Since there is only one sample in each class on the training set of Fa, the dictionary learning methods are not suitable to use in this case and the performance is too poor to be presented.

conducted based on deep learning. In the experiments, the Caffe deep learning framework [75] was used as the pre-processing to learn the deep features from the sample images. After the deep features were obtained, we further used the subspace learning methods (i.e. PCA, UDFS, SLDA, LPP, RFS, SpLpp, RR, LDA and the proposed JSLR) and dictionary learning methods (i.e. LCLE-DL and LC-KSVD2) to perform further feature extraction and then the nearest neighbor classifier was used for classification.

For AR and the standard subsets of the FERET databases, the dimension of extracted features based on the deep convolutional neural network (CNN) is 512 while that of the LFW database is 1024. For the standard FERET dataset, the Fa subset was used as the gallery set while the Fb, Fc, Dup1 and Dup2 were used as the probe sets. The LFW database contains images from 5,749 subjects in the uncontrolled environment, which makes it as a challenging recognition task. 158 subjects with total 4,324 images are selected from LFW-a subset and used in our experiment as the LFW-a subset is the aligned version of LFW database. The sample images on this database are shown in Fig. 3(a).

The experimental results on AR databases is listed in Table VIII. For the standard subsets of the FERET database, the best recognition rates corresponding different methods are

shown in Table IX, in which the accuracy on both the original dimensions (i.e. 512) and 216 dimensions (i.e. half of 512) are listed. The results in Table VIII and IX clearly show that the performance of the proposed JSLR is better than that of Deep-NN. This indicates that JSLR is able to extract discriminative information from deep features and further achieve higher recognition rate. The experimental results on LFW database are shown in Fig. 3(c) and Table X. In Fig. 3(c), the reason why no curves of LCLE-DL and LC-KSVD2 present is that no PCA is used as pre-processing to reduce the dimension of the input data to a specific value, the dictionary learning methods (i.e. LCLE-DL and LC-KSVD2) can only obtain recognition rate corresponding to the original dimension (i.e. 1024). Therefore, we cannot obtain the recognition rate curve versus the dimension variations for the two methods. From Fig. 3(c) and Table X, we can know that JSLR outperforms other compared methods again.

The convergence curves of the proposed JSLR on all databases are shown in Fig. 3(d) and Fig. 4. In these figures, the objective function value corresponding to each iteration is denoted as e^y where y is the values marked on the vertical coordinate. The convergence curves on all databases indicate that the proposed method can converge after several iterations.

TABLE X
THE PERFORMANCE (RECOGNITION RATE, STANDARD DEVIATION AND DIMENSION) OF ALL METHODS ON THE LFW FACE DATABASE BASED ON DEEP LEARNING

Training samples	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	Deep-NN	JSLR ($\alpha_1=0$)	JSLR
3	96.46	96.50	97.59	91.13	96.96	89.29	95.06	55.60	95.39	93.94	96.45	95.51	98.19
	± 0.18	± 0.18	± 0.23	± 0.67	± 0.19	± 1.02	± 0.43	± 3.36	± 0.57	± 0.96	± 0.16	± 0.39	± 0.20
	15*5	16*5	31*5	19*5	31*5	9*5	31*5	29*5	40*5	40*5	197*5	40*5	40*5
5	96.79	96.80	98.58	97.82	98.39	96.67	96.48	69.64	97.51	95.48	96.78	98.39	98.71
	± 0.23	± 0.22	± 0.22	± 0.20	± 0.26	± 0.24	± 0.32	± 2.22	± 0.32	± 0.30	± 0.20	± 0.14	± 0.20
	14*5	13*5	31*5	21*5	31*5	9*5	31*5	25*5	40*5	40*5	198*5	40*5	39*5
7	97.06	97.12	99.01	98.85	98.96	98.10	97.07	70.47	97.95	97.16	97.08	98.91	98.92
	± 0.11	± 0.15	± 0.14	± 0.17	± 0.10	± 0.20	± 0.32	± 2.71	± 0.35	± 0.38	± 0.12	± 0.14	± 0.13
	20*5	13*5	31*5	25*5	31*5	9*5	25*5	26*5	40*5	40*5	197*5	39*5	39*5

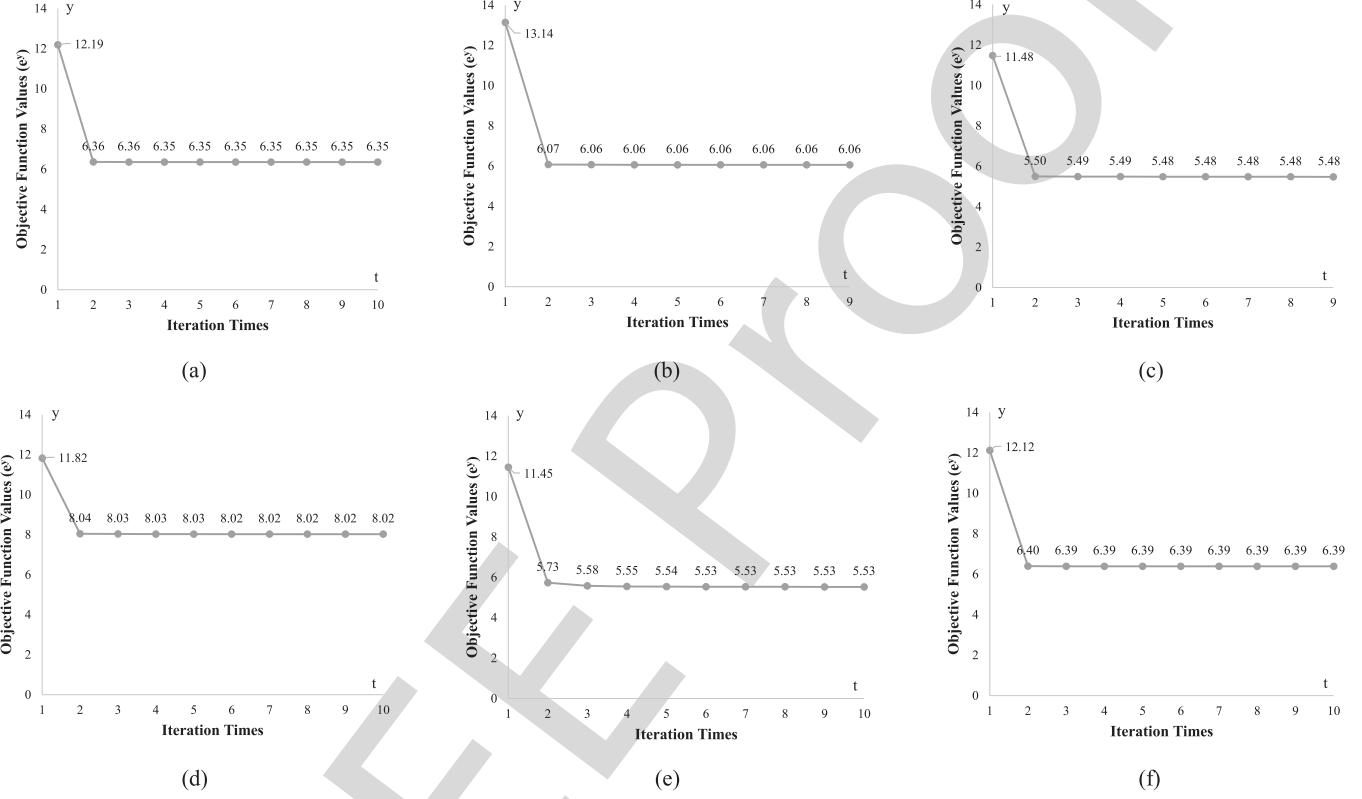


Fig. 4. Examples of the convergence curves of JSLR on (a) FERET, (b) AR, (c) CMU PIE, (d) USPS, (e) PIE67 \times 170 and (f) LFW database, respectively.

790 D. Experimental Results and Discussions

791 The comparison among the proposed JSLR, classical PCA,
792 RR, LDA, SLDA, LPP, $L_{2,1}$ -norm based methods
793 (UDFS, RFS) and dictionary learning methods (LCLE-DL, LC-
794 KSVD2) has been presented using recognition rates on these
795 databases: FERET, AR, CMU PIE, Yale and USPS. From the
796 results, we reveal the following interesting points:
797

- 798 1) In all experiments, including the face databases (FERET,
799 AR, CMU PIE, Yale) and non-facial database (USPS),
800 JSLR consistently achieves higher recognition rates than
801 other methods. These results are in line with the theoretical
802 analysis of JSLR that it obtains discriminative information
803 with joint sparsity and takes local geometric structure of
804 dataset into consideration to perform feature selection and
805 extraction.
806 2) JSLR is able to encode more discriminating information
807 in the low-dimensional face subspace since the local

808 geometric structure is considered to be more effective than
809 the global structure for feature extraction and feature se-
810 lection in some cases. The reason why JSLR outperforms
811 the local structure learning method such as LPP and SpLpp
812 is that JSLR utilizes α -norm regularization for feature se-
813 lection and feature extraction to obtain the discriminative
814 information for face recognition.

- 815 3) As it can be seen from the Fig. 2(d), (e) and (f), the tra-
816 ditional regression methods and/or their extensions can
817 obtain only c projections for feature extraction and classi-
818 fication, which is not enough to achieve high recognition
819 rates. Note that the number of classes in CMU PIE and
820 Yale is 68 and 15, respectively. Therefore the numbers of
821 projections obtained by LDA are 67 and 14, and the num-
822 bers of projections obtained by RR are 68 and 15 on CMU
823 PIE and Yale databases, respectively. Fig. 2(e) and (f) show
824 that the recognition rates of RR and LDA achieve their top

TABLE XI
THE COMPUTATIONAL COST (UNIT: S) OF DIFFERENT METHODS

Data set (l)	PCA	UDFS	SLDA	LPP	RFS	SpLPP	RR	LDA	LCLE-DL	LC-KSVD2	JSLR ($\alpha=0$)	JSLR
FERET ($l=4$)	0.0720	0.0390	14.7585	0.0201	1.7227	23.8032	0.1456	0.0642	1.7157	21.1985	0.0334	0.0301
AR ($l=4$)	0.0305	0.0566	7.2524	0.0097	0.5550	39.4835	0.1484	0.0607	1.1984	11.4715	0.0242	0.0237
CMU PIE ($l=4$)	0.0087	0.0767	7.6509	0.0375	0.1741	457.7060	0.0383	0.0483	0.4201	6.1862	0.0341	0.0906
Yale ($l=4$)	0.0019	0.0053	0.1162	0.0018	0.0161	0.4617	0.1251	0.0039	0.0505	1.1682	0.0028	0.0035
USPS ($l=400$)	0.0105	0.0433	4.0517	0.2033	131.6732	240.4646	0.0091	0.0376	6.7102	347.7577	0.1423	0.1462
PIE67 \times 170 ($l=80$)	0.1224	0.0738	168.2569	0.0289	285.8122	1.3378 $\times 10^3$	0.1190	0.2580	7.7944	630.7555	0.2123	0.2864
LFW ($l=4$)	0.0229	0.0661	48.6211	0.0157	1.4418	255.8663	0.0319	0.0984	1.1866	18.8623	0.1565	0.2012
Fa	0.0234	0.4945	1.3758 $\times 10^4$	0.0611	10.0917	3.8150 $\times 10^3$	0.0236	0.0371	18.1813	1.5484 $\times 10^3$	0.3679	0.2854

recognition rates using all the projections (we copy the final recognition rate to full fill all the dimensions listed on the horizontal axis). Thus the recognition rates no more increase after the number of dimension reaches 67 and 14 for LDA and 68 and 15 for RR on CMU PIE and Yale face databases, respectively. These figures show that the lack of enough projection of LDA and RR limits their performances. However, JSLR can break through this limitation and obtain more projections. This is the potential reason for JSLR to achieve higher recognition rates. In addition, the experimental result on USPS database with more than 4000 samples (as shown in Table VI) indicates the robustness and effectiveness of JSLR in dealing with large-sample size problem.

- 4) The $L_{2,1}$ -norm based methods such as JSLR, RFS and UDFS are robust to outliers in dataset and they guarantee the joint sparsity. However, JSLR obtains the best recognition rates when there are variations on lighting condition and face expressions. This indicates that JSLR is more robust than RFS and UDFS in feature extraction and selection when there exists variations on lighting condition and face expressions. In addition, the experimental results based on deep learning techniques presented in Table VIII and Table IX indicate the good performance of the proposed JSLR.
- 5) Experimental results indicate that the proposed JSLR performs better than the dictionary learning methods (LCLE-DL and LC-KSVD2). The reason is that JSLR guarantees the joint sparsity for discriminant feature selection or extraction in different cases. The Comparison between Deep-NN and Deep-JSLR shows that JSLR can further enhance the discriminative power of the deeply learned features based on CNNs in face recognition and character recognition tasks.
- 6) Table XI presents the computational time (unit: second) of each method on different data sets. From Table XI, we can know that the proposed JSLR based on $L_{2,1}$ -norm minimization is fast convergent and the computational cost is much less than the L_1 -norm based methods (i.e.

SLDA, SpLPP). The essential reason is that both SLDA and SpLPP use the least-angle regression method to compute the sparse solution and the iteration times are more than the proposed method. Moreover, the projections of SLDA and SpLPP are computed one by one while JSLR can simultaneously compute a set of jointly sparse projections. Thus, the proposed JSLR is efficient and effective for computer vision and pattern recognition.

- 7) From Tables III and IV, we can see that when more training samples are used, the recognition rates of all methods are higher than that when less training samples are used. However, it does not mean that more training samples can definitely help to obtain higher recognition rate. As shown in Table IV, when 20/24 training samples are used, the recognition rates of all methods become lower compared to the case when 19/24 training samples are used. The potential reason for this phenomenon is that too many training samples may lead to overfitting and thus all methods obtain poorer performance in the testing stage.

VII. CONCLUSION

Motivated by previous works that $L_{2,1}$ -norm regularization is able to obtain joint sparsity, and the local geometric information can enhance feature selection capability, in this paper, we propose a novel method called JSLR for feature extraction and selection. With $L_{2,1}$ -norm regularization and locality preserving property, JSLR can obtain any number of discriminative projections for feature selection, which addresses the drawback in LDA and ridge regression. Theoretical analyses show the close relationship of JSLR and ridge regression, which also guarantees the effectiveness of JSLR in feature extraction and selection. In order to obtain the optimal solution of JSLR, we propose an iterative algorithm which is proved to be convergent. In addition, the computational complexity of the algorithm is also presented. The performance of JSLR on several well-known face databases shows that it outperforms the classical principle component analysis methods, traditional sparse learning methods and recently proposed $L_{2,1}$ -norm regularization methods.

900 APPENDIX
901 PROOF OF THEOREM 1

902 From Eq. (10), we have

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{X}^T \mathbf{B} \mathbf{A}^T\|_F^2 \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{Y} - 2\mathbf{B}^T \mathbf{X} \mathbf{Y} \mathbf{A} + \mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B}). \end{aligned}$$

903 By setting the derivatives of the above problem with respect
904 to \mathbf{B} equaling to 0, we have

$$\mathbf{B} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}.$$

905 Let \mathbf{B}^* represents the optimal solution of Eq. (10), then

$$\mathbf{B}^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A}.$$

906 Since $\mathbf{P}^0 = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}$, we have $\mathbf{B}^* = \mathbf{P}^0 \mathbf{A}$. As matrix \mathbf{A} is a rotation matrix, then the subspace spanned by \mathbf{B}^* in Eq. (10) is the same as that spanned by \mathbf{P}^0 in Eq. (1), namely, $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$.

907 Suppose $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, by the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, we have $\mathbf{B}^* = \mathbf{U} \frac{1}{D} \mathbf{V}^T \mathbf{Y} \mathbf{A}$. Since the optimal solution of Eq. (3) is $\mathbf{P}^* = \mathbf{U} \frac{D}{D^2 + \alpha I} \mathbf{V}^T \mathbf{Y}$, then we can find that the subspaces spanned by \mathbf{B}^* and \mathbf{P}^* have the same base matrix \mathbf{U} and the only difference is that there is a weighted rotation matrix, which does not affect the spanned subspace. Thus, we can say $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$.

917 If $\alpha \rightarrow 0$, we have

$$\mathbf{P}^* \mathbf{A} = \mathbf{U} \frac{D}{D^2 + \alpha I} \mathbf{V}^T \mathbf{Y} \mathbf{A} \rightarrow \mathbf{U} \frac{1}{D} \mathbf{V}^T \mathbf{Y} \mathbf{A} = \mathbf{B}^*.$$

918 Thus, for any two pattern vectors \mathbf{x}_i and \mathbf{x}_j , since $\mathbf{A}^T \mathbf{A} = \mathbf{I}$,
919 the distance of the two points obtained by using the two sub-
920 spaces (i.e. \mathbf{B}^* and \mathbf{P}^*) for feature extraction is invariant. That
921 is,

$$\|(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{B}^*\|_2 = \|(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{P}^* \mathbf{A}\|_2 = \|(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{P}^*\|_2,$$

922 which indicates that the performance of using the two metric
923 matrices derived by \mathbf{B}^* and \mathbf{P}^* for classification will be the
924 same.

925 For Eq. (19), if $\alpha_1 \rightarrow 0$ and $\alpha_2 \rightarrow 0$, we have

$$\begin{aligned} \bar{\mathbf{B}} &= (\mathbf{X} \mathbf{X}^T + \alpha_1 \mathbf{D}_B + \alpha_2 \mathbf{X} (\bar{\mathbf{D}} - \bar{\mathbf{W}}) \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} \\ &\rightarrow (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y} \mathbf{A} = \mathbf{B}^*, \end{aligned}$$

926 namely, $\text{span}(\mathbf{B}^*) \rightarrow \text{span}(\bar{\mathbf{B}})$.

927 Since $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^0)$, $\text{span}(\mathbf{B}^*) = \text{span}(\mathbf{P}^*)$, then $\text{span}(\bar{\mathbf{B}}) \rightarrow \text{span}(\mathbf{P}^0)$, $\text{span}(\bar{\mathbf{B}}) \rightarrow \text{span}(\mathbf{P}^*)$.

929 REFERENCES

- [1] G. V. Lashkia and L. Anthony, "Relevant, irredundant feature selection and noisy example elimination," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 34, no. 2, pp. 888–897, Apr. 2004.
- [2] T. W. S. Chow, P. Wang, and E. W. M. Ma, "A new feature selection scheme using a data distribution factor for unsupervised nominal data," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 38, no. 2, pp. 499–509, Apr. 2008.
- [3] J. Qian, J. Yang, and Y. Xu, "Local structure-based image decomposition for feature extraction with applications to face recognition," *IEEE Trans. Image Process.*, vol. 22, no. 9, pp. 3591–3603, Sep. 2013.
- [4] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *J. Comput. Graph. Stat.*, vol. 15, no. 2, pp. 265–286, 2006.
- [5] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. fisherfaces: Recognition using class specific linear projection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 7, pp. 711–720, Jul. 1997.
- [6] W. Zheng and X. Tang, "A robust algorithm for generalized orthonormal discriminant vectors," in *Proc. IEEE 18th Int. Conf. Pattern Recognit.*, 2006, vol. 2, pp. 784–787.
- [7] G. Dai and Y. Qian, "A gabor direct fractional-step LDA algorithm for face recognition," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2004, vol. 1, pp. 61–64.
- [8] I. Dagher, "Incremental PCA-LDA algorithm," in *Proc. IEEE Int. Conf. Comput. Intell. Meas. Syst. Appl.*, 2010, pp. 97–101.
- [9] R. Raghavendra, B. Dorizzi, A. Rao, and G. H. Kumar, "Designing efficient fusion schemes for multimodal biometric systems using face and palmprint," *Pattern Recognit.*, vol. 44, no. 5, pp. 1076–1088, 2011.
- [10] M. Eskandari and Ö. Toygar, "Selection of optimized features and weights on face-Iris fusion using distance images," *Comput. Vision Image Understanding*, vol. 137, pp. 63–75, 2015.
- [11] F. Zhang, J. Yang, J. Qian, and Y. Xu, "Nuclear norm-based 2-DPCA for extracting features from images," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 10, pp. 2247–2260, Oct. 2015.
- [12] W. Wang, Y. Yan, F. Nie, S. Yan, and N. Sebe, "Flexible manifold learning with optimal graph for image and video representation," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2664–2675, Jun. 2018.
- [13] F. Nie, S. Yang, R. Zhang, and X. Li, "A general framework for auto-weighted feature selection via global redundancy minimization," *IEEE Trans. Image Process.*, vol. 28, no. 5, pp. 2428–2438, Dec. 2018.
- [14] H. Zhang, J. Yang, J. Xie, J. Qian, and B. Zhang, "Weighted sparse coding regularized nonconvex matrix regression for robust face recognition," *Inf. Sci.*, vol. 394, pp. 1–17, 2017.
- [15] X. Ning, W. Li, B. Tang, and H. He, "BULDP: Biomimetic uncorrelated locality discriminant projection for feature extraction in face recognition," *IEEE Trans. Image Process.*, vol. 27, no. 5, pp. 2575–2586, Feb. 2018.
- [16] R. Tibshirani, "Regression shrinkage and selection via the lasso: A retrospective," *J. Roy. Stat. Soc.: Ser. B (Stat. Method.)*, vol. 73, no. 3, pp. 273–282, 2011.
- [17] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. Roy. Stat. Soc.: Ser. B (Stat. Method.)*, vol. 67, no. 2, pp. 301–320, 2005.
- [18] I. T. Jolliffe, N. T. Trendafilov, and M. Uddin, "A modified principal component technique based on the lasso," *J. Comput. Graph. Stat.*, vol. 12, no. 3, pp. 531–547, 2003.
- [19] A. d'Aspremont, L. E. Ghaoui, M. I. Jordan, and G. R. Lanckriet, "A direct formulation for sparse PCA using semidefinite programming," in *Proc. Adv. Neural Inf. Process. Syst.*, 2005, pp. 41–48.
- [20] J. Feng, L. Jiao, F. Liu, T. Sun, and X. Zhang, "Unsupervised feature selection based on maximum information and minimum redundancy for hyperspectral images," *Pattern Recognit.*, vol. 51, pp. 295–309, 2016.
- [21] Z. Li and J. Tang, "Unsupervised feature selection via nonnegative spectral analysis and redundancy control," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5343–5355, Dec. 2015.
- [22] L. Clemmensen, T. Hastie, D. Witten, and B. Ersbøll, "Sparse discriminant analysis," *Technometrics*, vol. 53, no. 4, pp. 406–413, 2011.
- [23] Z. Qiao, L. Zhou, and J. Z. Huang, "Sparse linear discriminant analysis with applications to high dimensional low sample size data," *Int. J. Appl. Math.*, vol. 39, no. 1, 2009.
- [24] Z. Zeng, X. Wang, J. Zhang, and Q. Wu, "Semi-supervised feature selection based on local discriminative information," *Neurocomputing*, vol. 173, pp. 102–109, 2016.
- [25] C. Shi, Q. Ruan, G. An, and R. Zhao, "Hessian semi-supervised sparse feature selection based on $l_{2,1/2}$ -matrix norm," *IEEE Trans. Multimedia*, vol. 17, no. 1, pp. 16–28, Jan. 2014.
- [26] C. Ding and D. Tao, "Robust face recognition via multimodal deep face representation," *IEEE Trans. Multimedia*, vol. 17, no. 11, pp. 2049–2058, Nov. 2015.
- [27] A. Majumdar and R. K. Ward, "Classification via group sparsity promoting regularization," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, 2009, pp. 861–864.
- [28] A. Majumdar and R. K. Ward, "Robust classifiers for data reduced via random projections," *IEEE Trans. Syst., Man, Cybern., Part B Cybern.*, vol. 40, no. 5, pp. 1359–1371, Oct. 2010.
- [29] A. Majumdar and R. K. Ward, "Fast group sparse classification," *Can. J. Electr. Comput. Eng.*, vol. 34, no. 4, pp. 136–144, 2009.

- 1015 [30] X. He, D. Cai, S. Yan, and H.-J. Zhang, "Neighborhood preserving embedding," in *Proc. 10th IEEE Int. Conf. Comput. Vision Volume 1*, 2005, vol. 2, pp. 1208–1213.
- 1016 [31] D. Cai *et al.*, "Isometric projection," in *Proc. Assoc. Advancement Artif. Intell.*, 2007, pp. 528–533.
- 1017 [32] X. He and P. Niyogi, "Locality preserving projections," in *Proc. Advances Neural Inf. Process. Syst.*, 2004, pp. 153–160.
- 1018 [33] X. He, S. Yan, Y. Hu, P. Niyogi, and H.-J. Zhang, "Face recognition using laplacianfaces," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 3, pp. 328–340, Mar. 2005.
- 1019 [34] D. Cai, X. He, and J. Han, "Spectral regression: A unified approach for sparse subspace learning," in *Proc. IEEE Int. Conf. Data Mining*, 2007, pp. 73–82.
- 1020 [35] X. Liu, L. Wang, J. Zhang, J. Yin, and H. Liu, "Global and local structure preservation for feature selection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 6, pp. 1083–1095, Jun. 2014.
- 1021 [36] S. Liao *et al.*, "Discriminant analysis via joint euler transform and $\ell_{2,1}$ -norm," *IEEE Trans. Image Process.*, vol. 27, no. 11, pp. 5668–5682, Nov. 2018.
- 1022 [37] J. Yang *et al.*, "Nuclear norm based matrix regression with applications to face recognition with occlusion and illumination changes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 1, pp. 156–171, Jan. 2017.
- 1023 [38] J. Qian, L. Lei, Y. Jian, F. Zhang, and Z. Lin, "Robust nuclear norm regularized regression for face recognition with occlusion," *Pattern Recognit.*, vol. 48, no. 10, pp. 3145–3159, 2015.
- 1024 [39] L. Luo, J. Yang, J. Qian, Y. Tai, and G. F. Lu, "Robust image regression based on the extended matrix variate power exponential distribution of dependent noise," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 9, pp. 2168–2182, Sep. 2017.
- 1025 [40] L. Luo, J. Yang, J. Qian, and Y. Tai, "Nuclear- ℓ_1 norm joint regression for face reconstruction and recognition with mixed noise," *Pattern Recognit.*, vol. 48, no. 12, pp. 3811–3824, 2015.
- 1026 [41] J. Chen, J. Yang, L. Luo, J. Qian, and W. Xu, "Matrix variate distribution-induced sparse representation for robust image classification," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 10, pp. 2291–2300, Oct. 2015.
- 1027 [42] M. Yang, L. Zhang, J. Yang, and D. Zhang, "Regularized robust coding for face recognition," *IEEE Trans. Image Process.*, vol. 22, no. 5, pp. 1753–1766, May 2013.
- 1028 [43] Y. Li, J. Si, G. Zhou, S. Huang, and S. Chen, "FREL: A stable feature selection algorithm," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 7, pp. 1388–1402, Jul. 2015.
- 1029 [44] W. Yang, C. Sun, and W. Zheng, "A regularized least square based discriminative projections for feature extraction," *Neurocomputing*, vol. 175, pp. 198–205, 2016.
- 1030 [45] H. Pu and G. Gao, "Parameterless reconstructive discriminant analysis for feature extraction," *Neurocomputing*, vol. 190, pp. 50–59, 2016.
- 1031 [46] Y. Yang, Z.-J. Zha, Y. Gao, X. Zhu, and T.-S. Chua, "Exploiting web images for semantic video indexing via robust sample-specific loss," *IEEE Trans. Multimedia*, vol. 16, no. 6, pp. 1677–1689, Oct. 2014.
- 1032 [47] S. Xiang, F. Nie, G. Meng, C. Pan, and C. Zhang, "Discriminative least squares regression for multiclass classification and feature selection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 23, no. 11, pp. 1738–1754, Nov. 2012.
- 1033 [48] J. Han, Z. Sun, and H. Hao, " ℓ_0 -norm based structural sparse least square regression for feature selection," *Pattern Recognit.*, vol. 48, no. 12, pp. 3927–3940, 2015.
- 1034 [49] J. Yang and C. Ong, "An effective feature selection method via mutual information estimation," *IEEE Trans. Syst. Man Cybern. Part B Cybern.*, vol. 42, no. 6, pp. 1550–1559, Dec. 2012.
- 1035 [50] Z. Zeng, H. Zhang, R. Zhang, and C. Yin, "A novel feature selection method considering feature interaction," *Pattern Recognit.*, vol. 48, no. 8, pp. 2656–2666, 2015.
- 1036 [51] Y. Hui and Y. Jian, "Sparse discriminative feature selection," *Pattern Recognit.*, vol. 48, no. 5, pp. 1827–1835, 2015.
- 1037 [52] Y. Liu, F. Tang, and Z. Zeng, "Feature selection based on dependency margin," *IEEE Trans. Cybern.*, vol. 45, no. 6, pp. 1209–1221, Jun. 2015.
- 1038 [53] C. Hou, F. Nie, X. Li, D. Yi, and Y. Wu, "Joint embedding learning and sparse regression: A framework for unsupervised feature selection," *IEEE Trans. Cybern.*, vol. 44, no. 6, pp. 793–804, Jun. 2014.
- 1039 [54] F. Nie, H. Huang, C. Xiao, and C. H. Q. Ding, "Efficient and robust feature selection via joint $\ell_2, 1$ -norms minimization," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2010, pp. 1813–1821.
- 1040 [55] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*. Berlin/Heidelberg: Springer, Jul. 2003.
- 1041 [56] Y. Yang, H. T. Shen, Z. Ma, Z. Huang, and X. Zhou, "L_{2,1}-norm regularized discriminative feature selection for unsupervised learning," in *Proc. Int. Joint Conf. Artif. Intell.*, 2011, pp. 1589–1594.
- 1042 [57] Y. Yang, Z. Ma, Y. Yang, F. Nie, and H. T. Shen, "Multitask spectral clustering by exploring intertask correlation," *IEEE Trans. Cybern.*, vol. 45, no. 5, pp. 1083–1094, May 2015.
- 1043 [58] E. V. Den Berg and M. P. Friedlander, "Theoretical and empirical results for recovery from multiple measurements," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2516–2527, May 2010.
- 1044 [59] A. Majumdar and R. K. Ward, "Synthesis and analysis prior algorithms for joint-sparse recovery," in *Proc. IEEE Int. Conf. Acoust.*, 2012, pp. 3421–3424.
- 1045 [60] J. He, L. Ding, L. Jiang, and L. Ma, "Kernel ridge regression classification," in *Proc. Int. Joint Conf. Neural Netw.*, 2014, pp. 2263–2267.
- 1046 [61] Z. Zheng *et al.*, "Regression analysis of locality preserving projections via sparse penalty," *Inf. Sci.*, vol. 303, pp. 1–14, 2015.
- 1047 [62] G. Shikkenawis and S. K. Mitra, "On some variants of locality preserving projection," *Neurocomputing*, vol. 173, pp. 196–211, 2016.
- 1048 [63] F. Nie, X. Wang, M. I. Jordan, and H. Huang, "The constrained laplacian rank algorithm for graph-based clustering," in *Proc. 30th AAAI Conf. Artif. Intell.*, 2016, pp. 1969–1976.
- 1049 [64] H. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos, "Multilinear principal component analysis of tensor objects for recognition," in *Proc. Int. Conf. Pattern Recognit.*, 2006, pp. 776–779.
- 1050 [65] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vision*, 2014, pp. 818–833.
- 1051 [66] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv preprint arXiv:1409.1556*.
- 1052 [67] Q. Wang, Z. Qin, F. Nie, and Y. Yuan, "Convolutional 2D LDA for nonlinear dimensionality reduction," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 2929–2935.
- 1053 [68] Z. Jiang, L. Zhe, and L. S. Davis, "Learning a discriminative dictionary for sparse coding via label consistent K-SVD," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2011, pp. 1697–1704.
- 1054 [69] Z. Li, Z. Lai, Y. Xu, J. Yang, and D. Zhang, "A locality-constrained and label embedding dictionary learning algorithm for image classification," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 2, pp. 278–293, Feb. 2017.
- 1055 [70] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *Proc. Eur. Conf. Comput. Vision*, 2016, pp. 499–515.
- 1056 [71] P. J. Phillips, H. Moon, S. A. Rizvi, and P. J. Rauss, "The FERET evaluation methodology for face-recognition algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 10, pp. 1090–1104, Oct. 2000.
- 1057 [72] A. Martinez and R. Benavente, "The AR face database," CVC, West Lafayette, IN, USA, Tech. Rep. 24, 1998.
- 1058 [73] T. Sim, S. Baker, and M. Bsat, "The CMU pose, illumination and expression database of human faces," Carnegie Mellon Univ.: Pittsburgh, Pennsylvania, USA, Tech. Rep. CMU-RI-TR-OI-02, 2001.
- 1059 [74] G. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," Tech. Rep. 07-49, UMass, 2007.
- 1060 [75] Y. Jia *et al.*, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. 22nd ACM Int. Conf. Multimedia*, 2014, pp. 675–678.