EMBEDDED CLUSTERING VIA ROBUST ORTHOGONAL LEAST SQUARE DISCRIMINANT ANALYSIS

Rui Zhang, Feiping Nie^{*}, and Xuelong Li

School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL) Northwestern Polytechnical University, Xi'an 710072, Shaanxi, China.

ABSTRACT

In this paper, a novel embedded clustering (EC) method is derived from the perspective of extending the supervised orthogonal least square discriminant analysis (OLSDA) method to the unsupervised case, which proves to be closely related to k-means. To achieve more statistical and structural properties, the robust learning of unsupervised OLSDA is investigated to further derive the unsupervised robust OLSDA (ROLSDA) problem. For the convenience of solving the proposed ROLSDA problem, re-weighted counterpart of ROLSDA is utilized with self-adaptive weight, such that the smaller weight would be assigned to the term with larger outliers automatically. Consequently, aforementioned EC method is proposed with not only the robust outliers but also the optimal weighted cluster centroids. Comparative experiments are presented to show the effectiveness of the EC method under the proposed ROLSDA problem.

Index Terms— Embedded clustering, least square discriminant analysis, robust learning, re-weighted problem.

1. INTRODUCTION

Orthogonal least square discriminant analysis (OLSDA) [1] serves as a pretty significant supervised technique for dimensionality reduction [2, 3, 4, 5, 6, 7, 8]. OLSDA method is derived from minimizing the discriminative information in data. In other words, OLSDA method virtually minimizes the within-class scatter matrix from the viewpoint of the least square regression. To achieve great discriminative power, the projected data in the same class are expected to be regressed to a single vector with the purpose of seeking an orthogonal projection, such that the error of square regression is minimized.

In this paper, former supervised OLSDA method is further extended to the unsupervised case. By virtue of revisiting both OLSDA and orthogonal centroid method (OCM) [9], a brand new least square form of unsupervised OLSDA can be derived. Besides, the proposed unsupervised OLSDA problem is proved to be closely related to k-means [10]. Moreover, robust OLSDA (ROLSDA) problem is further proposed and investigated via robust learning[11, 12, 13, 14]. To unravel the proposed ROLSDA problem, re-weighted counterpart of ROLSDA is utilized with self-adaptive weight, such that the smaller weight would be assigned to the term with larger outliers automatically. Accordingly, a novel embedded clustering (EC) method is derived from the proposed ROLSDA problem with both the robust outliers and the optimal weighted cluster centroids.

The rest of the paper is organized as follows. Section 2 revisits the OLSDA and the OCM methods. Section 3 extends supervised OLSDA to the unsupervised OLSDA with a novel least square form. Consequently, EC method is derived from the unsupervised robust OLSDA problem. In Section 4, comparative experiments are presented to show the effectiveness of the proposed EC method. Section 5 concludes the paper.

2. ORTHOGONAL LEAST SQUARE DISCRIMINANT ANALYSIS REVISITED

Given the training dataset $\mathscr{X} = \{(\mathbf{x_i}, \mathbf{c_i}) | \mathbf{x_i} \in \mathbb{R}^{d \times 1}; i = 1, 2, \ldots, n\}$ and related data matrix $\mathbf{X} = [\mathbf{x_1}, \mathbf{x_2}, \ldots, \mathbf{x_n}] \in \mathbb{R}^{d \times n}$ which are distributed in **c** different classes with dimension **d** and data number **n**, each data point $\mathbf{x_i}$ is associated with a class label $\mathbf{c_i} \in \{1, 2, \ldots c\}$. Denote \mathscr{X}_i is the dataset of **i**-th class and $\mathbf{n_i}$ is the number of data points in **i**-th class, then the within-class scatter matrix $\mathbf{S_w}$, the between-class scatter matrix $\mathbf{S_b}$ and the total-class scatter matrix $\mathbf{S_t}$ are defined as follows:

$$\begin{cases} \mathbf{S}_{\mathbf{w}} = \sum_{i=1}^{c} \sum_{\mathbf{x} \in \mathscr{X}_{i}} (\mathbf{x} - \bar{\mathbf{x}}_{i}) (\mathbf{x} - \bar{\mathbf{x}}_{i})^{\mathbf{T}} \\ \mathbf{S}_{\mathbf{b}} = \sum_{i=1}^{c} \mathbf{n}_{i} (\bar{\mathbf{x}}_{i} - \bar{\mathbf{x}}) (\bar{\mathbf{x}}_{i} - \bar{\mathbf{x}})^{\mathbf{T}} \\ \mathbf{S}_{\mathbf{t}} = \sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{\mathbf{x}}) (\mathbf{x}_{i} - \bar{\mathbf{x}})^{\mathbf{T}} \end{cases}$$
(1)

where $\bar{\mathbf{x}}_i = \frac{1}{n_i} \sum_{\mathbf{x}_j \in \mathscr{X}_i} \mathbf{x}_j$ is the class-specific mean of the *i*-th class and $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ is the global mean. The least

^{*}Corresponding author. Email: feipingnie@gmail.com.

squared loss function could be illustrated as:

$$\varepsilon = \|\mathbf{T_1} - \mathbf{T_2}\|_{\mathbf{F}}^2 \tag{2}$$

where $\|\cdot\|_F^2$ represents Frobenius norm. Define $A^{(t)}=\frac{1}{n}\mathbf{1}\mathbf{1}^T$ and

$$\mathbf{A_{ij}^{(w)}} = \left\{ \begin{array}{ll} \frac{1}{\mathbf{n_{c_i}}} & \mathbf{c_i} = \mathbf{c_j} \\ \mathbf{0} & otherwise \end{array} \right.$$

where $\mathbf{1} = (1, 1, \dots, 1)^{\mathbf{T}} \in \mathbb{R}^{n \times 1}$. By substituting $\mathbf{T}_1 = \mathbf{W}^{\mathbf{T}}\mathbf{X}$ and $\mathbf{T}_2 = \mathbf{W}^{\mathbf{T}}\mathbf{X}\mathbf{A}^{(\mathbf{w})}$ in Eq. (2) with orthogonal constraint $\mathbf{W}^{\mathbf{T}}\mathbf{W} = \mathbf{I}, \mathbf{W} \in \mathbb{R}^{d \times k}$, we have

$$\begin{aligned} \|\mathbf{W}^{\mathrm{T}}\mathbf{X} - \mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{A}^{(\mathbf{w})}\|_{\mathrm{F}}^{2} &= \mathrm{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{X}(\mathbf{I} - \mathbf{A}^{(\mathbf{w})})^{2}\mathbf{X}^{\mathrm{T}}\mathbf{W}) \\ &= \mathrm{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{S}_{\mathbf{w}}\mathbf{W}) \end{aligned}$$

which is the objective function of OLSDA discussed in [1]. Similarly, we could further set $T_1 = W^T X$ and $T_2 = W^T X A^{(t)}$ in Eq. (2) as

$$\begin{split} \|\mathbf{W}^{\mathrm{T}}\mathbf{X} - \mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{A}^{(\mathrm{t})}\|_{\mathrm{F}}^{2} &= \mathrm{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{X}(\mathbf{I} - \mathbf{A}^{(\mathrm{t})})^{2}\mathbf{X}^{\mathrm{T}}\mathbf{W}) \\ &= \mathrm{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{S}_{\mathrm{t}}\mathbf{W}) \end{split}$$
(4)

with orthogonal constraint $\mathbf{W}^{T}\mathbf{W} = \mathbf{I}$.

In terms of Eqs. (3) and (4), the objective function of OCM [9] could be rewritten as

$$\begin{split} \mathbf{Tr}(\mathbf{W}^{T}\mathbf{S}_{\mathbf{b}}\mathbf{W}) &= \mathbf{Tr}(\mathbf{W}^{T}(\mathbf{S}_{\mathbf{t}} - \mathbf{S}_{\mathbf{w}})\mathbf{W}) \\ &= \mathbf{Tr}(\mathbf{W}^{T}\mathbf{X}(\mathbf{I} - \mathbf{A}^{(t)} - \mathbf{I} + \mathbf{A}^{(w)})\mathbf{X}^{T}\mathbf{W}) \\ &= \mathbf{Tr}(\mathbf{W}^{T}\mathbf{X}(\mathbf{A}^{(w)} - \mathbf{A}^{(t)})\mathbf{X}^{T}\mathbf{W}) \end{split} (5) \\ \text{due to the idempotent matrices } \mathbf{A}^{(t)} \text{ and } \mathbf{A}^{(w)}, \text{ i.e., } (\mathbf{A}^{(t)})^{2} = \mathbf{A}^{(t)} \text{ and } (\mathbf{A}^{(w)})^{2} = \mathbf{A}^{(w)}. \end{split}$$

3. ROBUST ORTHOGONAL LEAST SQUARE DISCRIMINANT ANALYSIS

Based on Eqs. (4) and (5), the total-class scatter ${\bf S_t}$ and the between-class scatter ${\bf S_b}$ in (1) could be further reformulated into

$$\begin{cases} \mathbf{S}_{t} = \mathbf{X}\mathbf{H}\mathbf{X}^{T} \\ \mathbf{S}_{b} = \mathbf{X}\mathbf{H}\mathbf{Y}(\mathbf{Y}^{T}\mathbf{Y})^{-1}\mathbf{Y}^{T}\mathbf{H}\mathbf{X}^{T} \end{cases}$$
(6)

where $\mathbf{Y} \in \mathbb{R}^{\mathbf{n} \times \mathbf{c}}$ is the binary label matrix and $\mathbf{H} = \mathbf{I} - \mathbf{A}^{(t)}$ is the centering matrix. Based on Eq. (6), OLSDA in (3) could be rewritten as

$$\min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I}} \operatorname{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{S}_{\mathbf{w}}\mathbf{W}) = \min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I}} \operatorname{Tr}(\mathbf{W}^{\mathrm{T}}(\mathbf{S}_{t} - \mathbf{S}_{b})\mathbf{W})$$

$$= \min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I}} \operatorname{Tr}(\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H}(\mathbf{I} - \mathbf{Y}(\mathbf{Y}^{\mathrm{T}}\mathbf{Y})^{-1}\mathbf{Y}^{\mathrm{T}})\mathbf{H}\mathbf{X}^{\mathrm{T}}\mathbf{W})$$

$$= \min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I}} \|\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H}(\mathbf{I} - \mathbf{Y}(\mathbf{Y}^{\mathrm{T}}\mathbf{Y})^{-1}\mathbf{Y}^{\mathrm{T}})\|_{\mathbf{F}}^{2}$$

$$(7)$$

which serves as a special case of problem (2) by letting $T_1 = W^T X H$ and $T_2 = W^T X H Y (Y^T Y)^{-1} Y^T$, respectively.

• Relationship between OLSDA in (7) and k-means.

The k-means problem [10] can be recapitulated as

$$\min_{\mathbf{F},\mathbf{G}\in\mathbf{ind}} \|\mathbf{T}-\mathbf{F}\mathbf{G}^{\mathbf{T}}\|_{\mathbf{F}}^{2}$$
(8)

where each column of matrix $\mathbf{F} \in \mathbb{R}^{\mathbf{k} \times \mathbf{c}}$ represents the cluster centroid and each row of the indicative matrix $\mathbf{G} \in \mathbb{R}^{\mathbf{n} \times \mathbf{c}}$ demonstrates the binary label. If the associated label of data $\mathbf{T} \in \mathbb{R}^{\mathbf{k} \times \mathbf{n}}$ is known, i.e., indicative matrix \mathbf{G} is fixed as binary label \mathbf{Y} in (8), the *k*-means problem degenerates to

$$\min_{\mathbf{F}} \|\mathbf{T} - \mathbf{F}\mathbf{Y}^{\mathbf{T}}\|_{\mathbf{F}}^{2} = \|\mathbf{T} - \mathbf{T}\mathbf{Y}(\mathbf{Y}^{\mathbf{T}}\mathbf{Y})^{-1}\mathbf{Y}^{\mathbf{T}}\|_{\mathbf{F}}^{2}.$$
 (9)

By further replacing the data \mathbf{T} with the centralized projected data $\mathbf{W}^{T}\mathbf{X}\mathbf{H} \in \mathbb{R}^{\mathbf{k}\times\mathbf{n}}$ in Eq. (9), we notice that the problem (9) is same as the problem (7). In other words, OLSDA in (7) is equivalent to *k*-means in (8) when $\mathbf{T} = \mathbf{W}^{T}\mathbf{X}\mathbf{H}$ and $\mathbf{G} = \mathbf{Y}$. Accordingly, OLSDA in (7) could be naturally extended to the unsupervised case as

$$\min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I},\mathbf{F},\mathbf{G}\in\mathrm{ind}} \|\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H}-\mathbf{F}\mathbf{G}^{\mathrm{T}}\|_{\mathbf{F}}^{2}.$$
 (10)

• Embedded clustering via Robust OLSDA.

Based on the unsupervised OLSDA in (10), robust OLSDA (ROLSDA) could be proposed as

$$\min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I},\mathbf{F},\mathbf{G}\in\mathrm{ind}} \|\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H}-\mathbf{F}\mathbf{G}^{\mathrm{T}}\|_{2,1}.$$
 (11)

Motivated by [11, 12], we could utilize the re-weighted counterpart of ROLSDA in (11) for the convenience as

$$\min_{\mathbf{W}^{T}\mathbf{W}=\mathbf{I},\mathbf{F},\mathbf{G}\in\mathbf{ind}} \| (\mathbf{W}^{T}\mathbf{X}\mathbf{H} - \mathbf{F}\mathbf{G}^{T})\mathbf{D}^{\frac{1}{2}} \|_{\mathbf{F}}^{2}$$

$$= \min_{\mathbf{W}^{T}\mathbf{W}=\mathbf{I},\mathbf{F},\mathbf{G}\in\mathbf{ind}} \sum_{i=1}^{n} \mathbf{D}_{ii} \| \mathbf{W}^{T}\mathbf{x}_{i}^{(\mathbf{H})} - \mathbf{F}\mathbf{g}_{i} \|_{2}^{2}$$

$$(12)$$

where diagonal weight matrix **D** is to be updated iteratively in the algorithm with $\mathbf{D}_{ii} \leftarrow \frac{1}{2 || \mathbf{W}^T \mathbf{x}_i^{(H)} - \mathbf{Fg}_i ||_2}$. Besides, $\mathbf{x}_i^{(H)}$ and \mathbf{g}_i are i-th columns of **XH** and \mathbf{G}^T , respectively. Actually, the proposed re-weighted ROLSDA in (12) is connected with self-adaptive weight **D** such that re-weighted ROLSDA not only has robust outliers but achieves the optimal weighted cluster centroids as well.

Since the centroid matrix \mathbf{F} is free from any constraint in (12), the extreme value condition w.r.t. \mathbf{F} could be inferred that

$$\begin{aligned} &\frac{\partial \| (\mathbf{W}^{\mathrm{T}} \mathbf{X} \mathbf{H} - \mathbf{F} \mathbf{G}^{\mathrm{T}}) \mathbf{D}^{\frac{1}{2}} \|_{\mathbf{F}}^{2}}{\partial \mathbf{F}} = \mathbf{0} \\ &\Rightarrow \frac{\partial \mathbf{T} \mathbf{r} (\mathbf{F} \mathbf{G}^{\mathrm{T}} \mathbf{D} \mathbf{G} \mathbf{F}^{\mathrm{T}} - \mathbf{2} \mathbf{W}^{\mathrm{T}} \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{G} \mathbf{F}^{\mathrm{T}})}{\partial \mathbf{F}} = \mathbf{0} \end{aligned}$$
(13)
$$&\Rightarrow \mathbf{F} = \mathbf{W}^{\mathrm{T}} \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{G} (\mathbf{G}^{\mathrm{T}} \mathbf{D} \mathbf{G})^{-1}. \end{aligned}$$

Error rate			
Dataset	AR	GT	USPS
k-means[10]	0.3975	0.4117	0.4530
RMKMC[16]	0.3771	0.3912	0.4228
OLSDA[1]	0.3826	0.4006	0.4219
EC(our)	0.3611	0.3773	0.4200
NMI			
Dataset	AR	GT	USPS
k-means[10]	0.6556	0.6942	0.6258
RMKMC[16]	0.6723	0.7250	0.6547
OLSDA[1]	0.6730	0.7334	0.6124
EC(our)	0.6882	0.7483	0.6609

 Table 1. Comparisons of error rate and NMI for 4 methods as k-means[10], RMKMC[16] and unsupervised OLSDA[1] over 3 different datasets are performed.

Consequently, re-weighted ROLSDA in (12) degenerates to

$$\min_{\mathbf{W}^{\mathrm{T}}\mathbf{W}=\mathbf{I},\mathbf{G}\in\mathbf{ind}} \| (\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H} - \mathbf{F}\mathbf{G}^{\mathrm{T}})\mathbf{D}^{\frac{1}{2}} \|_{\mathbf{F}}^{2}$$
(14)

where $\mathbf{F} = \mathbf{W}^{T} \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{G} (\mathbf{G}^{T} \mathbf{D} \mathbf{G})^{-1}$. Via the coordinate blocking method, i.e., fixing each variable alternatively, reweighted ROLSDA in (14) could be solved correspondingly.

Case 1 (fixing G): The Lagrangian function of the problem (14) could be written as

$$\mathscr{L}(\mathbf{W}, \mathbf{\Lambda}) = \|(\mathbf{W}^{\mathrm{T}}\mathbf{X}\mathbf{H} - \mathbf{F}\mathbf{G}^{\mathrm{T}})\mathbf{D}^{\frac{1}{2}}\|_{\mathbf{F}}^{2} - \mathbf{Tr}(\mathbf{\Lambda}(\mathbf{W}^{\mathrm{T}}\mathbf{W} - \mathbf{I})).$$
(15)

Thus, the KKT condition can be derived as

$$\begin{aligned} \frac{\partial \mathscr{L}(\mathbf{W}, \mathbf{\Lambda})}{\partial \mathbf{W}} &= \mathbf{0} \\ \Rightarrow \mathbf{2}\mathbf{X}\mathbf{H}\mathbf{D}\mathbf{H}\mathbf{X}^{\mathrm{T}}\mathbf{W} - \mathbf{2}\mathbf{X}\mathbf{H}\mathbf{D}\mathbf{G}\mathbf{F}^{\mathrm{T}} &= \mathbf{2}\mathbf{W}\mathbf{\Lambda} \qquad (16) \\ \Rightarrow \mathbf{W}^{\mathrm{T}}(\mathbf{S}_{\mathbf{t}}^{(\mathrm{D})} - \mathbf{S}_{\mathbf{b}}^{(\mathrm{D})})\mathbf{W} &= \mathbf{\Lambda} \end{aligned}$$

where

$$\begin{cases} \mathbf{S}_{\mathbf{t}}^{(\mathbf{D})} = \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{H} \mathbf{X}^{\mathbf{T}} \\ \mathbf{S}_{\mathbf{b}}^{(\mathbf{D})} = \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{G} (\mathbf{G}^{\mathbf{T}} \mathbf{D} \mathbf{G})^{-1} \mathbf{G}^{\mathbf{T}} \mathbf{D} \mathbf{H} \mathbf{X}^{\mathbf{T}} \end{cases}$$
(17)

Eqs. (14) and (16) in whole imply that \mathbf{W} is the matrix of eigenvector corresponding to the first \mathbf{k} smallest eigenvalues of $\mathbf{S}_{\mathbf{t}}^{(\mathbf{D})} - \mathbf{S}_{\mathbf{b}}^{(\mathbf{D})}$ defined in (17).

Case 2 (fixing W): To obtain the indicative matrix **G**, re-weighted ROLSDA in (14) could be further derived as

$$\min_{\mathbf{G} \in \mathbf{ind}} \| (\mathbf{W}^{T} \mathbf{X} \mathbf{H} - \mathbf{F} \mathbf{G}^{T}) \mathbf{D}^{\frac{1}{2}} \|_{\mathbf{F}}^{2}$$

$$= \min_{\mathbf{g}_{i} \in \mathbf{ind}} \sum_{i=1}^{n} \mathbf{D}_{ii} \| \mathbf{W}^{T} \mathbf{x}_{i}^{(\mathbf{H})} - \mathbf{F} \mathbf{g}_{i} \|_{2}^{2}$$

$$\geq \sum_{i=1}^{n} \min_{\mathbf{g}_{i} \in \mathbf{ind}} \mathbf{D}_{ii} \| \mathbf{W}^{T} \mathbf{x}_{i}^{(\mathbf{H})} - \mathbf{F} \mathbf{g}_{i} \|_{2}^{2}$$

$$(18)$$

Input: data matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$ and cluster number c. **Output:** indicative matrix $\mathbf{G} \in \mathbb{R}^{n \times c}$.

1 Initialize weight matrix $\mathbf{D} = \mathbf{I}$ with random initial guess $\mathbf{G} \in \mathbf{ind}$ and random orthogonal subspace $\mathbf{W} \in \mathbb{R}^{\mathbf{d} \times \mathbf{k}}$, i.e., $\mathbf{W}^{T}\mathbf{W} = \mathbf{I}$;

2 while not converge do

- 3 Update $\mathbf{F} \leftarrow \mathbf{W}^{\mathbf{T}} \mathbf{X} \mathbf{H} \mathbf{D} \mathbf{G} (\mathbf{G}^{\mathbf{T}} \mathbf{D} \mathbf{G})^{-1}$;
- $\begin{array}{c|c|c} \mathbf{4} & \quad \mathbf{for} \ \mathbf{i} = \mathbf{1} : \mathbf{n} \ \mathbf{do} \\ \mathbf{5} & \quad Update \\ & \quad \mathbf{g}_i \leftarrow \arg\min_{\mathbf{g}_i \in \mathbf{ind}, \mathbf{1_c}^{\mathrm{T}} \mathbf{g}_i = \mathbf{1}} \| \mathbf{W}^{\mathrm{T}} \mathbf{x}_i^{(\mathrm{H})} \mathbf{F} \mathbf{g}_i \|_2^2; \\ \mathbf{6} & \quad \mathbf{end} \\ \mathbf{7} & \quad \mathbf{for} \ \mathbf{i} = \mathbf{1} : \mathbf{n} \ \mathbf{do} \\ \mathbf{8} & \quad | \quad Update \ \mathbf{D}_{ii} \leftarrow \frac{1}{2 \| \mathbf{W}^{\mathrm{T}} \mathbf{x}_i^{(\mathrm{H})} \mathbf{F} \mathbf{g}_i \|_2}; \\ \mathbf{9} & \quad \mathbf{end} \\ \mathbf{1} & \quad \mathbf{1} \leftarrow \mathbf{0}^{(\mathrm{D})} \quad \mathbf{MUDMWT} \\ \end{array}$
- 10 Update $\mathbf{S}_{\mathbf{t}}^{(\mathbf{D})} \leftarrow \mathbf{XHDHX^{T}};$ 11 Update
- 12 $\mathbf{S}_{\mathbf{b}}^{(\mathbf{D})} \leftarrow \mathbf{X}\mathbf{H}\mathbf{D}\mathbf{G}(\mathbf{G}^{T}\mathbf{D}\mathbf{G})^{-1}\mathbf{G}^{T}\mathbf{D}\mathbf{H}\mathbf{X}^{T};$ 12 Update $\mathbf{W} \leftarrow \operatorname*{arg\,min}_{\mathbf{W}^{T}\mathbf{W}=\mathbf{I}}\mathbf{Tr}(\mathbf{W}^{T}(\mathbf{S}_{\mathbf{t}}^{(\mathbf{D})} - \mathbf{S}_{\mathbf{b}}^{(\mathbf{D})})\mathbf{W});$

13 end

14 return G;

ş

Algorithm 1: Embedded clustering (EC) method under the proposed ROLSDA problem in (11).

which indicates that indicative vector $\mathbf{g}_i \in \mathbb{R}^{c \times 1}$ of i-th data point could be determined by individually solving

$$\min_{\mathbf{g}_i \in \mathbf{ind}} \| \mathbf{W}^{\mathrm{T}} \mathbf{x}_i^{(\mathrm{H})} - \mathbf{F} \mathbf{g}_i \|_2^2 \text{ s.t. } \mathbf{1_c}^{\mathrm{T}} \mathbf{g}_i = 1$$
(19)

where $\mathbf{1_c} = (\mathbf{1}, \dots, \mathbf{1})^{\mathbf{T}} \in \mathbb{R}^{\mathbf{c} \times \mathbf{1}}$. Based on Eq. (14), Case 1 and 2, embedded clustering (EC) method could be summarized in Algorithm 1 to solve the proposed ROLSDA in (11).

4. EXPERIMENT

We utilize 9 benchmark datasets as AT&T, $COIL_{20}$, $COIL_{100}$, FEI, IMM, YALE, AR, GT and USPS in the experimental part to further compare the clustering accuracy and NMI under different reduced dimensionality. Besides, clustering accuracy and NMI are defined as following.

• **Clustering accuracy.** The clustering accuracy in the experiment can be computed as

$$\mathbf{Accuracy} = \frac{1}{n} \max(\sum_{\mathbf{R_k}, \mathbf{O_m}} \mathbf{M}(\mathbf{R_k}, \mathbf{O_m}))$$

where $\mathbf{R_k}$ stands for the k-th cluster in the final result and $\mathbf{O_m}$ stands for the true m-th class. $\mathbf{M}(\mathbf{R_k}, \mathbf{O_m})$ stands for the number of entities, which originally serve as the input data



Fig. 1. The clustering accuracy comparisons are performed for the PCA[15]+k-means[10] method, the PCA[15]+RMKMC[16] method, the unsupervised OLSDA[1] method and the proposed EC method with two baselines as the k-means[10] method and the RMKMC[16] method under 6 benchmark datasets. (a) AT&T. (b) COIL₂₀. (c) COIL₁₀₀. (d) FEI. (e) IMM. (f) YALE.

of the m-th class and are assigned to the k-th cluster in the final result.

• Normalized mutual information. The normalized mutual Information (NMI) serves as an index to determine the quality of the clusters. The NMI index is defined as

$$\mathbf{NMI} = \frac{\sum_{k_1=1}^{c} \sum_{k_2=1}^{c} n_{\mathbf{C_{k_1}} \cap \mathbf{C_{k_2}}} \log(\frac{nn_{\mathbf{C_{k_1}} \cap \mathbf{C_{k_2}}}{n_{k_1} n_{k_2}})}{\sqrt{\sum_{k_1=1}^{c} n_{k_1} \log{(\frac{n_{k_1}}{n})}} \sqrt{\sqrt{\sum_{k_2=1}^{c} n_{k_2} \log{(\frac{n_{k_2}}{n})}}}$$

where $\mathbf{n_{k_1}}, (1 \leq \mathbf{k_1} \leq \mathbf{c})$ denotes the number of data in cluster $\mathbf{C_{k_1}}$ and $\mathbf{n_{k_2}}, (1 \leq \mathbf{k_2} \leq \mathbf{c})$ denotes the number of data in cluster $\mathbf{C_{k_2}}$ with $\mathbf{n_{C_{k_1} \cap C_{k_2}}}$ being the number of data in the intersection set $\mathbf{C_{k_1}} \cap \mathbf{C_{k_2}}$.

In Fig. 1, first six datasets of the selected datasets mentioned above are utilized to compare the clustering accuracy of four approaches as PCA[15] + k-means[10], PCA[15] + RMKMC[16], unsupervised OLSDA[1] and the proposed EC under different reduced dimensionality with k-means[10] and RMKMC[16] being the baselines. In particular, PCA[15] + k-means[10] and PCA[15] + RMKMC[16] denote that lowdimensional data $W^T X$ is achieved in advance via PCA[15] for further clustering via the k-means[10] and RMKMC[16] methods, respectively.

In Tab. 1, average least error rates are recorded for unsupervised OLSDA and the proposed EC method to compare with the error rates obtained by the k-means and RMKMC methods under the rest of the datasets. Besides, the associated **NMI** index is also recorded in Tab. 1.

According to the results in Fig. 1 and Tab. 1, we could conclude that:

1. From Fig. 1, the clustering accuracy of the proposed EC method is better than those of other approaches including PCA + k-means, PCA + RMKMC, unsupervised OLSDA and two baselines under most reduced dimensionality.

2. From Tab. 1, the proposed EC method is consistently better than *k*-means, RMKMC and unsupervised OLSDA on both error rate and NMI index.

5. CONCLUSION

In this paper, a novel embedded clustering method is proposed via extending the supervised orthogonal least square discriminant analysis method to the robust unsupervised case, which is closely related to k-means. Besides, re-weighted counterpart of unsupervised robust orthogonal least square discriminant analysis problem is utilized for the convenient optimization with self-adaptive weight. Consequently, the embedded clustering method is proposed with robust outliers and optimal weighted cluster centroids. The effectiveness and the superiority of the proposed embedded clustering method are further verified both empirically and analytically.

6. REFERENCES

- F. Nie, S. Xiang, Y. Liu, C. Hou, and C. Zhang, "Orthogonal vs. uncorrelated least squares discriminant analysis for feature extraction," *Pattern Recognition Letters*, vol. 33, no. 5, pp. 485–491, 2012.
- [2] F. Nie, H. Huang, C. Ding, D. Luo, and H. Wang, "Robust principal component analysis with non-greedy 11norm maximization," in *Proceedings*. International Joint Conference on Artificial Intelligence, 2011, pp. 1433– 1438.
- [3] F. Nie, D. Xu, W. Tsang, and C. Zhang, "Flexible manifold embedding: a framework for semi-supervised and unsupervised dimension reduction," *IEEE Transactions* on *Image Processing*, vol. 19, no. 7, pp. 1921–1932, 2010.
- [4] T. Zhou and D. Tao, "Double shrinking sparse dimension reduction," *IEEE Transactions on Image Processing*, vol. 22, no. 1, pp. 244–257, 2013.
- [5] X. Li, Y. Pang, and Y. Yuan, "11-norm based 2dpca," *IEEE Trans. on System, Man and Cybernetics, Part B*, vol. 40, no. 4, pp. 1170–1175, 2010.
- [6] Y. Pang, D. Tao, Y. Yuan, and X. Li, "Binary twodimensional pca," *IEEE Trans. on System, Man and Cybernetics, Part B*, vol. 38, no. 4, pp. 1176–1180, 2008.
- [7] C. Zhang, F. Nie, and S. Xiang, "A general kernelization framework for learning algorithms based on kernel pca," *Neurocomputing*, vol. 73, no. 4-6, pp. 959–967, 2010.
- [8] T. Luo, C. Hou, D. Yi, and J. Zhang, "Discriminative orthogonal elastic preserving projections for classification," *Neurocomputing*, vol. 179, pp. 54–58, 2016.
- [9] H. Park, M. Jeon, and J. B. Rosen, "Lower dimensional representation of text data based on centroids and least squares," *BIT Numer. Math.*, vol. 43, no. 2, pp. 427–448, 2003.
- [10] J. Macqueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings*. Berkeley Symposium on Mathematical Statistics and Probability, 1967, vol. 1, pp. 281–297.
- [11] F. Nie, J. Yuan, and H. Huang, "Optimal mean robust principal component analysis," in *Proceedings*. International Conference on Machine Learning, 2014, pp. 2755–2763.
- [12] W. Jiang, F. Nie, and H. Huang, "Robust dictionary learning with capped 11 -norm," in *Proceedings*. International Joint Conference on Artificial Intelligence, 2015, pp. 3590–3596.

- [13] G. Lan, C. Hou, and D. Yi, "Robust feature selection via simultaneous capped l2-norm and l2;1-norm minimization," in *IEEE*. International Conference on Big Data Analysis, 2016.
- [14] M. Qian and C. Zhai, "Robust unsupervised feature selection," in *Proceedings*. International Joint Conference on Artificial Intelligence, 2013, pp. 1621–1627.
- [15] I. T. Jolliffe, *Principal Component Analysis*, Springer-Verlag, New York, NY, USA, 2nd edition, 2002.
- [16] X. Cai, F. Nie, and H. Huang, "Multi-view k-means clustering on big data," in *Proceedings*. International Joint Conference on Artificial Intelligence, 2013, pp. 2598–2604.