

Facial Expression Recognition Based on LLENet

Dan Meng¹, Guitao Cao^{*,1,2}

¹ School of Computer Science and
Software Engineering

East China Normal University
Shanghai, China 200062

Corresponding Author Email:

*gtcao@sei.ecnu.edu.cn

Zhihai He² IEEE Fellow

²Department of Electrical and
Computer Engineering

University of Missouri
Columbia, USA MO65211

Wenming Cao^{+,3}

³College of Information Engineering
Shenzhen University

Shenzhen, China 518060

Co-corresponding Author Email:

⁺wmcao@szu.edu.cn

Abstract—Facial expression recognition plays an important role in lie detection, and computer-aided diagnosis. Many deep learning facial expression feature extraction methods have a great improvement in recognition accuracy and robustness than traditional feature extraction methods. However, most of current deep learning methods need special parameter tuning and ad hoc fine-tuning tricks. This paper proposes a novel feature extraction model called Locally Linear Embedding Network (LLENet) for facial expression recognition. The proposed LLENet first reconstructs image sets for the cropped images. Unlike previous deep convolutional neural networks that initialized convolutional kernels randomly, we learn multi-stage kernels from reconstructed image sets directly in a supervised way. Also, we create an improved LLE to select kernels, from which we can obtain the most representative feature maps. Furthermore, to better measure the contribution of these kernels, a new distance based on kernel Euclidean is proposed. After the procedure of multi-scale feature analysis, feature representations are finally sent into a linear classifier. Experimental results on facial expression datasets (CK+) show that the proposed model can capture most representative features and thus improves previous results.

Index Terms—Face expression recognition, deep learning, kernel distance, locally linear embedding (LLE).

I. INTRODUCTION

For many tasks in pattern recognition, such as hand-written digit recognition, face recognition and object recognition, discovering suitable representations of images is a critical problem. Among these visual tasks, facial expression recognition is considered to be the most challenging and long-standing topic since it is widely used in areas of lie detection, surveillance, identity authentication, access control and human-computer interaction. Classical feature extraction methods for face images include histogram of oriented gradients (HOG), local binary pattern (LBP), Gabor features and their fusions. These low-level feature extraction methods often lack the ability of capturing the most useful information directly from data samples. As a result, although these feature representations work well on certain face recognition tasks, they can hardly cope with uncontrolled environments (e.g. occlusion, large variations in illumination and pose, and accessory changes).

Over the past few years, feature extraction methods based on manifold learning and deep learning have been made breakthrough on a wide range of vision tasks [1], [2], [3]. It is worth notice that, Chan [3] developed a very simple DNN for

computer vision classification tasks, called PCANet. PCANet [3] utilized the most basic and easy processing components to emulate DNN, which is able to train and adapt to various tasks easily. However, such a simple architecture ignores the valuable known label information of the sample images during training stage, thus affects the recognition accuracy. Besides, Chan [3] applied linear dimension reduction method (PCA) to find kernels (detectors) with maximum responses, which might be unsuitable for facial expression recognition, given face images' nonlinear manifold structure [1]. Enlightened by PCANet [3], in this paper we introduce a novel feature extraction model called LLENet for face expression recognition to take advantage of features learned through deep learning method while overcome drawbacks of PCANet [3]. The architecture of the proposed feature extraction model is simple, since it only consists of a procedure of image reconstruction sets (IRS), a multi-stage feature convolutional layer and a multi-scale feature analysis layer. Besides, we remove back propagation from standard DNN, so we do not need special parameter tuning and ad hoc fine-tuning tricks.

The rest of the paper is organized as follows. Section II presents the proposed LLENet, and experimental results are provided in Section III followed by conclusions in Section IV.

II. OUR METHOD (LLENET)

Given face video sequences or images, we use a facial detector to localize the face region of each frame or image. The localized face region is then cropped and normalized to a fixed height and width, generating a face image. To extract feature representations, face images in IRS are convoluted with learned kernels before multi-scale feature analysis. These feature representations are then fed into classification model. Figure 1 summarizes our proposed LLENet.

A. Image Reconstruction Sets

Let $\{(X_i, L_i) | i = 1, 2, \dots, N_{train}\}$ be training images in the dataset, where $X_i \in R^{m \times n}$ represents the i -th training sample, and $L_i = 1, 2, \dots, C$ is the class label of X_i . For each training sample X_i , we take a $k_1 \times k_2$ patch with a step of 1 pixel around each pixel, and collect overlapped patches of i -th image $\tilde{X}_i = [x_{(i,1)}, x_{(i,2)}, \dots, x_{(i,\tilde{m}\tilde{n})}]$, where $x_{(i,j)}$

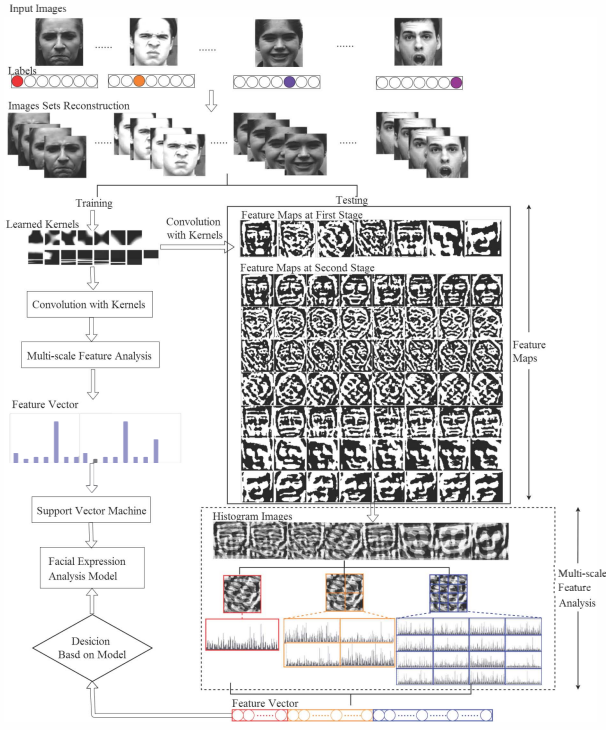


Fig. 1. Flowchart of face recognition framework based on LLENet.

denotes the j -th vectorized patch in X_i , $\tilde{m} = m - \lceil k_1/2 \rceil$, $\tilde{n} = n - \lceil k_2/2 \rceil$, and $\lceil z \rceil$ gives the smallest integer greater than or equal to z . Let $\tilde{X}_i^{(p)} \in R^{1 \times \tilde{m}\tilde{n}}$ demotes p -th row in \tilde{X}_i , which can be seen as a copy of the original image with some rows and columns removed from X_i when we reshape vector $\tilde{X}_i^{(p)}$ to a $\tilde{m} \times \tilde{n}$ matrix. We call $\tilde{X}_i = [\tilde{X}_i^{(1)}; \tilde{X}_i^{(2)}; \dots; \tilde{X}_i^{(k_1 k_2)}]$ a reconstructed image set, and the corresponding labels can be represented as $\tilde{L}_i = \{\tilde{L}_i^{(p)} = L_i, p \in [1, k_1 k_2]\}$. By constructing the same image set for all training samples, we finish reconstructing image sets $\tilde{X} = \{\tilde{X}_1; \tilde{X}_2; \dots; \tilde{X}_{N_{train}}\}$. The labels of the reconstructed image sets are denoted by $\tilde{L} = \{\tilde{L}_1; \tilde{L}_2; \dots; \tilde{L}_{N_{train}}\}$.

B. Improved LLE Based on the Modified Kernel Distance

Among all the image distance metrics d used in LLE algorithm, Euclidean distance d_E or its kernel form d_{Ek} are the most commonly used due to their simplicity. However, this distance measurement suffers from a high sensitivity even to small deformation. This paper proposes a modified kernel distance, which is called Kernel Nonlinear Distance d_{KND} . The modified distance d_{KND} is defined as:

$$d_{KND} = \left(\exp \frac{d_{Ek}^2}{\text{mean}(d_{Ek})} \right)^{\frac{1}{2}} + \alpha \max(d_{Ek}). \quad (1)$$

where $d_{Ek}(K_{IiCDM}(i, j)) = \sqrt{\langle \phi(K_{IiCDM}^i) - \phi(K_{IiCDM}^j), \phi(K_{IiCDM}^i) - \phi(K_{IiCDM}^j) \rangle}$ denotes kernel Euclidean distance, and we use classical

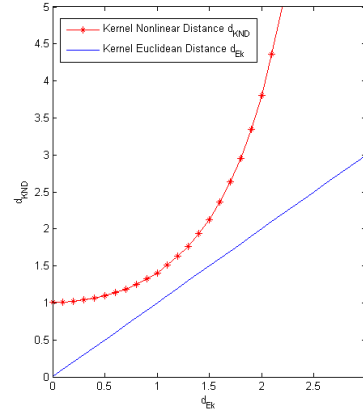


Fig. 2. Comparison of distance curve between d_{KND} and d_{Ek} , where $k_1 k_2 = 256$, and $\alpha = 1/256$.

Gauss kernel function as the kernel function. Besides, $K_{IiCDM} = \frac{S_{inter}}{S_{intra}}$, where S_{intra} and S_{inter} are the intra and inter scatter matrix for IRS images respectively. In Equation 1, $\text{mean}(d_{Ek})$ and $\max(d_{Ek})$ are the average and maximum kernel Euclidean distance between all pairs of data points respectively.

In order to have a better understanding of the advantage of the proposed nonlinear distance, Figure 2 shows a comparison of distance curves between d_{KND} and d_{Ek} when the original kernel Euclidean distance d_{Ek} increased linearly in the range of $[0, 3]$. From Figure 2, we can see that the modified distance d_{KND} brings distance between different convolutional kernels into an exponential growth mode. The reconstruction error of the original LLE can be rewritten as follows:

$$\varepsilon(W_i) = W_i^T d_{KND} W_i. \quad (2)$$

Thus objective cost function can be written as

$$\min J(V) = \text{tr}(V^T \Phi(K_{IiCDM}) M \Phi(K_{IiCDM})^T V). \quad (3)$$

The solution is known as the eigenvectors corresponding to the smallest d eigenvalues of M . Once V is obtained, for test samples, the nonlinear feature maps are given as $\tilde{Y}_i = \tilde{X}_i * V^T$.

C. Multi-scale features analysis

For each set of feature maps, we convert the feature maps belonging to the corresponding filter for the last stage into one histogram image whose every pixel is an integer in the range $[0, 255]$, and treated as a distinct “word”. For each set of B^τ output images X_{il}^τ belongs to the V_l^τ at final stage τ , we can obtain by:

$$X_{il}^\tau = X_i^{\tau-1} * \text{mat}_{k_1 k_2}(V_l^\tau)^T \quad (4)$$

where $i \in [1, N_{train} \prod_{t=1}^{\tau-1} B^t]$, $l \in [1, B^\tau]$, V_l^τ is the l -th filter at final stage τ , and $\text{mat}_{k_1 k_2}(V_l^\tau)$ is a function that

maps $V_l^\tau \in R^{k_1 k_2}$ to a maxtrix of dimension of $R^{k_1 \times k_2}$. Then histogram image $Hist_l^\tau$ can be expressed as:

$$Hist_{il} = \sum_{i=1}^{\prod_{t=1}^{\tau-1} B^t} (2^{\tau-1} \bmod 256) \times H(X_{il}^\tau). \quad (5)$$

where $H(X_{il})$ is a heviside step function. For every B^τ histogram image $Hist_{il}$, we construct a sequence of grids at resolutions $0, 1, \dots, Q$, such that the grid at level q has 2^q cells along each dimension, with a total of $G = \sum_{q=0}^Q 2^q$ cells. Let $Ch_g^q(Hist_{il})$ denotes the vector containing numbers of points from $Ch_g^q(Hist_{il})$ that fall into g -th cell at resolution q according to different “words”. We cascade $Ch_g^q(Hist_{il})$ to build a multi-scale feature map for histogram image as:

$$Ch(Hist_{il}) = [Ch_1^0(Hist_{il}), Ch_2^1(Hist_{il}), \dots, Ch_G^Q(Hist_{il})]. \quad (6)$$

where $Ch(Hist_{il}) \in R^{G \times 2^8}$. So feature representations for X_i can be expressed as:

$$f_{train}^i \doteq [Ch(Hist_{i1}), Ch(Hist_{i2}), \dots, Ch(Hist_{iB^\tau})]. \quad (7)$$

III. EXPERIMENT

To verify the proposed method, we performe experiments on CK+ to evaluate the performance of the proposed algorithm using linear SVM for classification. The CK+ [4] dataset consists of 593 sequences from 123 subjects of both sexes. We cropped facial expression images for each sequence by a face detector. Then we constructed three image pairs as done in [2] for each sequence. The first image in each pair corresponds to the onset image (neutral) and the other two images correspond to the two peak expression faces. As for the output facial expression images, the peak expression faces in image pairs were then normalized to 128×128 pixels, with labels corresponding to the peak expression faces.

We randomly select 70% samples to learn the convolutional kernels, and the rest samples extract features based on the learned kernels. Table I shows the comparison of the proposed LLENet and state-of-the-art methods. From Table I we can see that our proposed LLENet outperforms all the other methods and achieves better performance than the least STRBM [2]. Additionally, we present the confusion matrix of our method in Figure 3. From Figure, one can observe that LLENet achieves best performance on happy, disgust, and contempt, and worst performance on fear.

IV. CONCLUSION

In this paper, we proposed a novel feature extraction model based on supervised kernel learning and improved LLE for facial expression recognition. The proposed feature extraction model has three major components: image sets reconstruction (IRS), multi-stage supervised kernel learning, and multi-scale feature analysis. Experimental results show that the proposed feature extraction model is able to obtain robust and efficient feature representations for facial expression recognition task.

TABLE I
COMPARISON WITH STATE-OF-ART ON CK+ DATASET.

Method	Recognition Rate(%)
SPTS+CAPP [4]	83.30
CLM-SRI [5]	88.60
EAI [6]	82.60
LDN [7]	89.30
STM-ExpLet [8]	94.19
STRBM [2]	95.66
PCANet [3]	94.31
LLENet	96.94

	anger	contempt	disgust	fear	happy	sadness	surprise
anger	98.52	0.00	1.48	0.00	0.00	0.00	0.00
contempt	0.00	100.00	0.00	0.00	0.00	0.00	0.00
disgust	0.00	0.00	100.00	0.00	0.00	0.00	0.00
fear	0.00	0.00	1.33	85.33	5.33	0.00	0.00
happy	0.00	0.00	0.00	0.00	100.00	0.00	0.00
sadness	0.00	0.00	0.00	0.00	0.00	96.43	3.57
surprise	0.00	0.00	0.00	0.00	1.21	0.00	98.79

Fig. 3. Confusion Matrix of LLENet-2 on CK+ dataset.

ACKNOWLEDGMENT

This work was supported by Natural Science Foundation of China (61375015).

REFERENCES

- [1] G. Q. Wang, N. F. Shi, and Y. X. Shu, “Embedded manifold-based kernel fisher discriminant analysis for face recognition,” *Neural Process Letters*, vol. 43, no. 1, pp. 1–16, 2016.
- [2] S. Elaiwat, M. Bennamoun, and F. Boussaid, “A spatio-temporal RBM-based model for facial expression recognition,” *Pattern Recognition*, vol. 49, pp. 152–161, 2016.
- [3] T. H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma, “PCANet: A simple deep learning baseline for image classification,” *IEEE Transactions on Image Processing (TIP)*, vol. 24, no. 12, pp. 5017–5032, 2015.
- [4] P. Lucey, J. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, “Comprehensive data base (CK+): A complete dataset for action unit and emotion-specified expression,” in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshop*, Jun. 2010, pp. 94–101.
- [5] L. Jeni, D. Takacs, and A. Lorincz, “High quality facial expression recognition in video streams using shape related information only,” in *IEEE International Conference on Computer Vision (ICCV) Workshops*, Nov. 2011, pp. 2168–2174.
- [6] S. Yang and B. Bhanu, “Understanding discrete facial expressions in video using an emotion avatar image,” *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 42, no. 4, pp. 980–992, Aug. 2012.
- [7] A. R. Rivera, J. R. Castillo, and O. O. Chae, “Local directional number pattern for face analysis: Face and expression recognition,” *IEEE Transactions on Image Processing (TIP)*, vol. 22, no. 5, pp. 1740–1752, 2013.
- [8] M. Liu, S. Shan, R. Wang, and X. Chen, “Learning expression let son spatio-temporal manifold for dynamic facial expression recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.