

Orthogonal LDA in PCA Transformed Subspace

M. Mahadeva Prasad
Dept. of Electronics
University of Mysore
Post Graduate Centre
Hassan, INDIA.
prasada9@gmail.com

M. Sukumar
Dept. of Instr. Technology
S. J. College of
Engineering
Mysore, INDIA.
m_sukumar@rediffmail.com

A. G. Ramakrishnan
Dept. Of Electrical
Engineering
Indian Institute of Science
Bangalore, INDIA.
ramkiag@iisc.ernet.in

Abstract—The paper addresses the effectiveness of orthogonal linear discriminant analysis (OLDA) in a principal component analysis (PCA) transformed subspace. The performance of the technique is studied for writer independent recognition of online handwritten Kannada numerals. Experiments show that the performance of LDA and OLDA are better in a PCA transformed subspace compared to that of the original feature space. In addition, the recognition accuracies of the system with OLDA are marginally better than that of LDA in both the original feature space and the PCA transformed subspace. An average recognition accuracy of 96.9% is achieved on a database collected from 69 writers. To our knowledge, this is the first ever reported work on recognition of online handwritten Kannada numerals.

Keywords—Kannada numeral; online handwriting; OLDA;

I. INTRODUCTION

Form filling is one of the applications of online handwritten character recognition. Since the data for most form filling applications are collected in the field, there is a need for a portable and affordable system for data collection and also a robust handwritten data recognition system. Many of the form filling applications have fields exclusively meant for numeral data entry, such as date of birth, phone number, house number and PIN code. So, a recognition engine trained only on numerals will certainly perform better in such fields in terms of accuracy and recognition time than a system designed both for alphabets and numerals. Most Indian languages have their own symbols for the numerals and hence, the written data collected for form filling applications contain these numerals.

Attempts have been made to recognize offline handwritten numerals of Indian scripts like Hindi [1], Bangla [2,3] and Tamil [4]. In [5], offline handwritten numeral recognition for six popular Indian scripts is reported. In [4, 6, 7], attempts have also been made to recognize offline handwritten Kannada numerals. However, to our knowledge, no attempt has been made to recognize online handwritten numerals of Indian scripts. This motivated us to design an online handwritten Kannada numeral recognition system.

II. OLDA IN PCA TRANSFORMED SUBSPACE

Before Dimensionality reduction tools like PCA and LDA [8] have been successfully applied in many classification problems like face recognition, data mining and character recognition. In LDA, maximal class discrimination is achieved by finding the optimal discriminant vectors by maximizing the Fisher criterion function, i.e., the generalized Raleigh quotient. If S_w and S_b are within and between scatter matrices, the discriminant vectors derived by LDA are usually linearly dependent, since the matrix $S_w^{-1}S_b$ is asymmetric. Also for high-dimensional and small sample size problems, the traditional LDA faces the problem of singularity [9].

To eliminate the problem of dependency, the Foley–Sammon discriminant analysis (FSD) [10] was proposed, thereby improving the performance of LDA. But the theory of FSD and its computation are complex. Orthogonalized Fisher discriminant analysis (OFD) was proposed in [11] to overcome this complexity. In OFD, the orthogonal discriminant vectors are obtained by orthogonalizing the derived Fisher discriminant vectors. Similarly, to deal with the problem of singularity, null space LDA [12] and subspace LDA [9] have been proposed. In null space LDA, the discriminant vectors are calculated in the null space of the within-class scatter matrix and in the subspace LDA, the dimensionality of the original data is reduced before applying classical LDA. In [9], LDA is applied in the PCA transformed subspace to overcome the problem of singularity. Present work aims to study the performance of LDA and OLDA in PCA transformed subspace in the recognition of online handwritten Kannada numerals.

A. Foley–Sammon Discriminant Analysis

Suppose there are c known pattern classes. Let S_b and S_w be the between-class and within-class scatter matrices respectively. Let w_1, w_2, \dots, w_k be the set of discriminant vectors. LDA tries to determine the

transformation matrix $W = [w_1, w_2, \dots, w_k]$ by maximizing the Fisher criterion function:

$$J(W) = \frac{|W^T S_b W|}{|W^T S_w W|}. \quad (1)$$

Due to the asymmetric property of the matrix $S_w^{-1} S_b$, the derived Fisher discriminant vectors are usually linearly dependent. It has been proved that by eliminating this dependency, the performance of traditional LDA can be improved. Based on this fact, FSD tries to find a set of orthogonal discriminant vectors to constitute the transformation matrix.

Considering that first k Fisher vectors w_1, w_2, \dots, w_k ($k \geq 1$) have been obtained, the $(k+1)^{th}$ FSD vector w_{k+1} is obtained by maximizing the Rayleigh quotient with the following orthogonality constraints:

$$w_{k+1}^T w_i = 0, \quad i = 1, 2, \dots, k \quad (2)$$

It has been proved that the $(k+1)^{th}$ FSD vector w_{k+1} is the eigenvector of the matrix $S_w^{-1} P_k S_b$ corresponding to the largest Eigen value. The matrix P_k is calculated by,

$$P_k = I - [w_1, \dots, w_k] ([w_1, \dots, w_k]^T S_w^{-1} [w_1, \dots, w_k])^{-1} \times [w_1, \dots, w_k]^T S_w^{-1}. \quad (3)$$

For detailed analysis of FSD and its complexity, readers are referred to [10].

B. Orthogonal Fisher Discriminant Analysis

The OLDA was proposed to overcome the computational complexity of FSD. It proposes an easier way of obtaining the set of orthogonal discriminant vectors. Let w_1, w_2, \dots, w_r be r Fisher's discriminant vectors. Assume that the first k OFD vectors v_1, v_2, \dots, v_k ($1 \leq k \leq r-1$) have been obtained. The $(k+1)^{th}$ OFD vector can be obtained by the following relation:

$$v_{k+1} = w_{k+1} - \sum_{i=1}^k \frac{v_i^T w_{k+1}}{v_i^T v_i} v_i \quad (4)$$

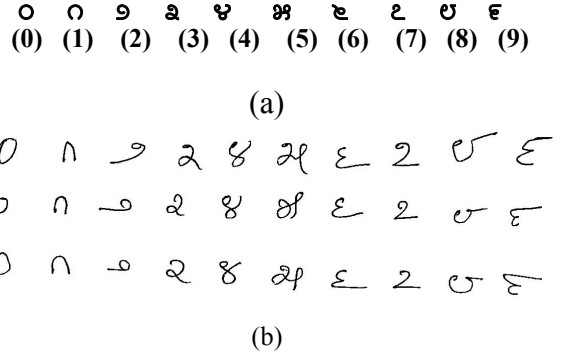


Figure 1. Kannada Numeral Set (a) Machine printed and (b) handwritten. Corresponding Indo-Arabic numerals are also shown for reference.

III. KANNADA SCRIPT

Kannada is the official language of Karnataka, one of the South Indian states of India. The script has 50 basic characters and 10 numerals. The complexity of the Kannada script and the challenges in designing a recognition engine for online handwritten Kannada characters can be found in [13]. Fig. 1 shows few samples of handwritten Kannada numerals along with a sample of machine printed numerals for reference.

IV. DATA COLLECTION AND PREPROCESSING

To our knowledge, there is no standard online handwritten Kannada numeral database available to carry out our experiments. So, we have created a Kannada numeral database. Tablet PC and G-Note 7000 digital pad are used to collect the data. Sixty nine native writers are requested to write the numbers in their own writing style without exercising any restrictions. Writers have written from one trial to a maximum of five trials. Totally there are one hundred and nine samples in each class.

The raw data of numerals consists of the x - and y -coordinate values corresponding to the pen-tip movements. The raw data is subjected to noise removal, re-sampling and normalization. Smoothing is performed by a moving average filter of size three. The smoothed data is re-sampled in space along the arc length by linear interpolation so that each numeral data consists of 30 points. Let the smoothed and sampled numeral data be represented by the sequence:

$$P = [p_1, p_2, \dots, p_{30}] \quad (5)$$

where the vector $p_i = (x_i, y_i)^T$ and x_i, y_i denote the horizontal and vertical coordinates. The sampled data are shifted and size normalized to get a new sequence:

$$Q = [q_1, q_2, \dots, q_{30}] \quad (6)$$

where the vector $q_i = (a_i, b_i)^T$ is given by,

$$a_i = (x_i - x_{\min}) / (x_{\max} - x_{\min}) \quad (7)$$

$$b_i = (y_i - y_{\min}) / (y_{\max} - y_{\min}) \quad (8)$$

where, $(x_{\min}, y_{\min})^T$ and $(x_{\max}, y_{\max})^T$ denote the minimum and the maximum horizontal and vertical coordinate values for the numeral under consideration.

V. FEATURE EXTRACTION

The shape of Kannada numerals motivated us to extract trajectory and deviation features from the size normalized sequence, Q .

A. Trajectory Features

The trajectory feature consists of distance and angle of each of the re-sampled points of the numeral from the origin of the normalized bounding box. Fig. 2 shows the way the trajectory features are extracted. A vector is defined from the origin to the starting point (a_1, b_1) of the normalized character. As we move along the character trace, the magnitude d_i of the vector with respect to the origin and the angle θ_i made by the vector with respect to the x -axis constitute the trajectory feature. A normalized distance feature vector D and an angle feature vector θ are computed from the sequence Q as follows.

$$D = \{d_1, d_2, d_3, \dots, d_{30}\} \quad (9)$$

$$\theta = \{\theta_1, \theta_2, \theta_3, \dots, \theta_{30}\} \quad (10)$$

where,

$$d_i = \frac{\sqrt{a_i^2 + b_i^2}}{\max(\sqrt{a_i^2 + b_i^2})} \quad \text{and} \quad \theta_i = \frac{\arg(a_i + jb_i)}{\max(\arg(a_i + jb_i))}$$

B. Deviation Features

The x and y deviations of normalized sample points of a character from its centroid are calculated using the following equations and are used as features.

$$X_d = (a_{d1}, a_{d2}, a_{d3}, \dots, a_{d30}) \quad (11)$$

$$Y_d = (b_{d1}, b_{d2}, b_{d3}, \dots, b_{d30}) \quad (12)$$

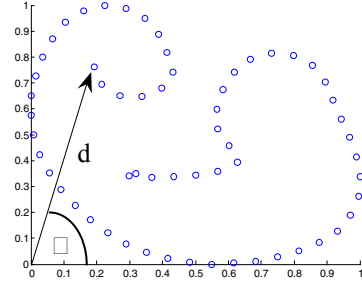


Figure 2. Extraction of trajectory features from a normalized Kannada character.

where,

$$a_{di} = (a_i - \mu_a) / \sigma_a, \quad b_{di} = (b_i - \mu_b) / \sigma_b$$

$$\mu_a = (1/30) \sum_{i=1}^{30} a_i, \quad \mu_b = (1/30) \sum_{i=1}^{30} b_i$$

$$\sigma_a = \sqrt{(1/29) \sum_{i=1}^{30} (\mu_a - a_i)^2} \quad \text{and} \quad \sigma_b = \sqrt{(1/29) \sum_{i=1}^{30} (\mu_b - b_i)^2}$$

VI. EXPERIMENTS AND RESULTS

Experiments are carried out in a writer independent mode. The collected database is segmented into two sets, belonging to different, disjoint sets of writers. For training, 70% of the database is used and the remaining 30% is used for testing.

The performance of OLDA and LDA are studied in the original feature space R^d and the PCA transformed subspace R^m ($m < d$). In order to retain the valuable discriminatory features in the R^m subspace, eigen vectors are selected such that $m = \text{rank}(S_t)$, where, S_t is the total scatter matrix. The recognition accuracies obtained by the experiments are tabulated in Table 1. The class-wise recognition accuracies of OLDA in PCA transformed subspace are given in Table 2. The results show that the performances of the recognition engines in the PCA transformed subspace are higher than those in the original feature space. The results also reveal that the performance of OLDA is marginally better than that of LDA.

Recognition of offline handwritten Kannada numerals has been carried out by different groups. Average recognition accuracies of 93%, 98.7% and 98.5% have been reported in [4], [5] and [6]. Since there is no previous work carried out to recognize online handwritten Kannada numerals, we could not compare our results with any other work.

VII. CONCLUSION

The performance of LDA and OLDA are studied in both original feature space and PCA transformed subspace. Experiments are carried out on writer independent recognition of online handwritten Kannada numerals. A maximum recognition accuracy of 96.9% is achieved. As part of our future work, we plan to improve the recognition accuracy by analyzing the confusion matrix. Further, experimentation of the writer dependent case and adaptability of the recognition system for any single writer will be investigated.

ACKNOWLEDGMENT

The authors thank Technology Development for Indian Languages (TDIL) Programme, Department of Information Technology, Ministry of Communication and Information Technology, Government of India for funding the collection of the Kannada online handwritten data used in this work. Thanks are also due to UGC for giving Teacher Fellowship to the first author of the paper to carry out research work.

REFERENCES

- [1] M. Hanmandlu, J. Grover, V. K. Madasu and S. Vasikarla, "Input fuzzy for the recognition of handwritten Hindi numeral", International Conference on Information Technology, 2007, vol. 2, pp. 208-213.
- [2] U. Battacharya, T. K. Das, A. Datta, S. K. Puri, and B. B. Choudhuri, "Recognition of handprinted Bangla numerals using neural network models", Lecture Notes in Computer Science, Vol. 2275/2002, pp. 139-161, Springer Berlin, 2002.
- [3] M. M. Honque, Islam Mohammad, M. M. Ali, "An efficient fuzzy method for Bangla handwritten numeral recognition", International conference on Electrical and Computer Engineering ICECE06, 2006, pp. 197-200.
- [4] S. V. Rajashekararadhya, P. Vanaja Ranjan, V. N. Manjunath Aradhya, "Isolated handwritten Kannada and Tamil numeral recognition: a novel approach", First International Conference on Emerging Trends in Engineering and Technology, ICETET 2008, pp. 1192-95.
- [5] U. Pal, N. Sharma, T. Wakabayashi, F. Kimura, "Hand-written numeral recognition of six popular Indian scripts", Ninth International conference on Document Analysis and Recognition ICDAR07, vol. 2, pp.749-753.
- [6] N. Sharma, U. Pal, F. Kimura, "Recognition of hand-written Kannada numerals", Proc. 9th International Conf. on Information Technology, 2006, pp. 133-136.
- [7] G. G. Rajaput and Mallikarjun Hangarge, "Recognition of isolated handwritten Kannada numerals based on image fusion method", PReMI07, 2007, LNCS. 4815, pp. 153-160.

TABLE I. Comparison of the recognition performance of LDA and OLDA on original and PCA features (Training set: 76 samples/class; Test set: 33 samples/class)

Tool	Original Feature Space		PCA Transformed Subspace	
	LDA	OLDA	LDA	OLDA
Recognition Accuracy %	92.4	92.8	95.9	96.9

TABLE II. Classwise recognition accuracy of OLDA in PCA transformed subspace

Digits	Writer Independent Recognition Accuracy (%)
೦	100.0
೧	100.0
೨	100.0
೩	100.0
೪	96.6
೫	96.6
೬	86.2
೭	96.6
೮	100.0
೯	93.1
Mean Accuracy	96.9

- [8] R. O. Duda, P. E. Hart, and D. Stork, "Pattern Classification", Second Edition, John Wiley and Sons (Asia) Pvt. Ltd., 2006.
- [9] P. N. Belhumeur, J. P. Hespanha, and David J. Kriegman, "Eigenfaces vs. Fisherfaces: recognition using class specific linear projection", IEEE Trans. Pattern Anal. Mach. Intell. 19 (7) (1997) 711-720
- [10] H. Foley and J. W. Sammon. "An optimal set of vectors", IEEE Transactions on Computers, 1975, vol 24, pp. 281-289.
- [11] Fengxi Song, Shuhai Liu, Jingyu Yang, "Orthogonalized Fisher discriminant", Pattern Recognition, 2005, Vol. 38, pp. 311 - 313.
- [12] L. F. Chen, H. Y. M. Liao, M. T. Ko, J. C. Lin, and G. J. Yu. "A new LDA-based face recognition system which can solve the small sample size problem". Pattern Recognition, 2000, Vol. 33, pp.1713-1726.
- [13] M. Mahadeva Prasad, M. Sukumar, A. G. Ramakrishnan, "Divide and conquer technique in online handwritten Kannada character recognition", Int. Workshop on Multilingual Optical Character Recognition Systems (MOCR), 2009, Barcelona, Spain.