

# Multiple Feature Extraction Methods for Handwritten Marathi Numeral Recognition

S. M. Mali

Department of Computer Science, MAEER' S Arts, Commerce and Science College, Pune

## Abstract

In this Paper, we present robust and novel methods for recognition of isolated handwritten Marathi numeral recognition. Density, Moments, Fourier Descriptor and chain codes are used as features. The recognition results are studied Support Vector Machine (SVM) classifiers. The system is experimented with our database of 12690 samples of Marathi handwritten numerals using fivefold cross validation method. Comparative study of different methods is presented.

## Keywords

Handwritten numerals, density, moment, Fourier Descriptor, Chain code, SVM classifier.

## 1. Introduction

Handwritten Marathi numeral recognition is a difficult problem, not only because of the great amount of variations in human handwriting, but also because certain numerals when handwritten look similar and sometimes they cannot be recognized manually also. These numerals written in Devanagari script are also used in many different Indian languages. Hence, as compared to other Indian scripts, the Devanagari script numerals find more applications in many areas including postal zip code processing and automatic data entry in forms in different languages

Many systems have been proposed for recognition of characters, both printed and handwritten, for Japanese, Persian and Arabic scripts, including Indic scripts [1, 2, 5, 8, 9, 10, 11, 13, and 15]. Some efficient system proposed for recognition of Devnagari includes [15, 12, 3, 6, 16, 17, 4]. We also found some work in Kannada [18, 7, 4].

## 2. Marathi numerals and pre-processing

Marathi is spoken by about 71 million people in Indian state of Maharashtra and neighboring states. Percentage wise it is the fourth widely spoken Indian language. Marathi script is written with Devanagari alphabet. Since the standard database for Marathi handwritten numerals is not available, to our

knowledge, the database is created with respect to variety in handwriting style. Data collection is done on a sheet specially designed for data collection. Data is collected from persons of different professions. A sample data sheet is shown in figure-1. Since, data is collected in a predefined format slant correction is assumed to be performed. The collected documents are scanned using scanner at 300 dpi to obtain gray scale images. The images are then binarized using Otsu's method. The speckles in the binarized images are removed using morphological erode and dilate operations. The numerals are cropped by fitting a minimum bounding box on the numeral. To bring uniformity among the numerals the cropped numeral is normalized to a size of 40x40 pixels. A total of 2500 binary digital images representing Marathi handwritten numerals are obtained. Each image represents a numeral (binary 1) that is unconstrained, isolated and clearly discriminated from the background (binary 0).

०	०	०	०	०	०	०	०	०	०	०
१	१	१	१	१	१	१	१	१	१	१
२	२	२	२	२	२	२	२	२	२	२
३	३	३	३	३	३	३	३	३	३	३
४	४	४	४	४	४	४	४	४	४	४
५	५	५	५	५	५	५	५	५	५	५
६	६	६	६	६	६	६	६	६	६	६
७	७	७	७	७	७	७	७	७	७	७
८	८	८	८	८	८	८	८	८	८	८
९	९	९	९	९	९	९	९	९	९	९

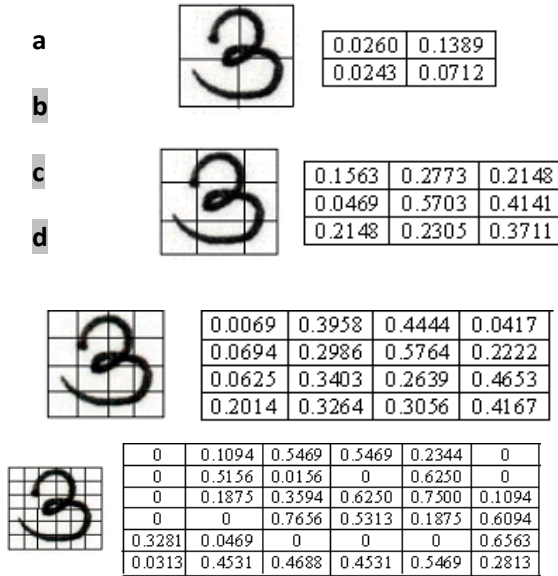
Fig. 1: Sample handwritten datasheet

## 3. Feature Extraction

Feature extraction facilitates extracting potential information about the numeral image. The information gathered forms the basis for feature vector that will be used in pattern classification. Following are different methods used for feature extraction.

### 3.1 Density

In the proposed method, the binary image representing the handwritten numeral is pre-processed as described in the Section 2 and is normalized to a size of 48 x 48 pixels. The size-normalized image is divided into n equal number of zones in sequence for n= 4,9,16 and 36 respectively.



**Figure 2 Zoning. Image divided into (a) 4 zones (b) 9 zones (c) 16 zones and (d) 36 zones An entry in a cell represents the density for that zone**

The density of the zone is computed by taking the ratio of total number of object pixels (i.e. pixels representing the numeral viz. binary 1) to total number of pixels in the zone (Equation 1). This is carried out for all the zones in the image. Totally 65 features are extracted from the image (Refer algorithm for details). An example is shown in Figure 2.

$$den(Z) = \frac{\text{No. of object pixels in the zone } Z}{\text{Total number of pixels in this zone}} \quad (1)$$

### 3.2 Moments

A lot of useful information about a binary object can be gained from the moments of the object. Most often moments are computed for connected components of binary images for the analysis of shape.

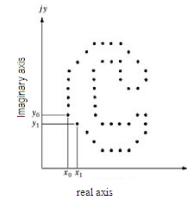
Moments may be computed for several orders. The central moments of third and Fourth order can be obtained using the following equation.

$$\mu_{pq} = \sum_x \sum_y (x - \bar{x})^p (y - \bar{y})^q f(x, y) \quad \text{for } p, q = 0, 1, 2, 3$$

such that  $p + q \leq 3$

### 3.3 Fourier Descriptors

The Fourier transform of a boundary representation (chain code, signature, complex boundary function) is used to represent the region's shape. One can low pass filter the boundary function spectrum without destroying the characteristic shape of the region. This means that only the amplitudes and phases of the low frequency components in the spectrum (i.e. the low-order Fourier coefficients) are required to characterize the basic shape of the object and they can be used as shape descriptors. Before calculating the Fourier descriptors input image must be segmented and boundary of the object must be determined. The boundary will be presented as an array of complex numbers which correspond to the pixels of the object boundary if the image is placed in the complex plane. Fourier descriptors are now calculated by combining Fourier transform coefficients of the complex array.



**Figure 3 A digital boundary and its representation as a complex sequence**

The points  $(x_0, y_0)$  and  $(x_1, y_1)$  are the first two points in the sequence. As shown in Figure 4.1, the K-point digital boundary in the xy-plane starting at an arbitrary point  $(x_0, y_0)$  co-ordinate pairs  $(x_0, y_0), (x_1, y_1), (x_2, y_2), \dots, (x_{k-1}, y_{k-1})$  are encountered in traversing the boundary say, in the counterclockwise direction. These co-ordinates can be expressed in the form of  $x(k) = x_k$  and  $y(k) = y_k$ . With this notation, the boundary itself can be represented as sequence of co-ordinates.

### 3.4 Chain Codes

Chain codes are one of the shape representations which are used to represent a boundary by a connected sequence of straight line segments of specified length

and direction. This representation is based on 4-connectivity or 8-connectivity of the segments [39]. The direction of each segment is coded by using a numbering scheme as shown in Figure below. The re-sampled boundary can be represented by a 4- or 8-code. The starting point can be arbitrarily chosen at the topmost, leftmost point of the boundary (Figure. 5.2)

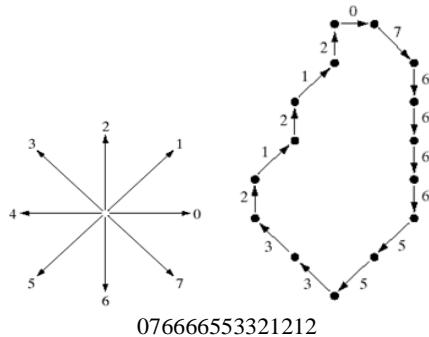


Figure 3-Directional Chain Code of re-sampled boundary

The chain code of the boundary can be normalized with respect to starting point by simply treating the chain code as a circular sequence of direction numbers and redefine the starting point so that the resulting sequence of numbers forms an integer of minimum magnitude. Further, to normalize the code for rotation (in angles that are integer multiples of directions) by using first difference of the chain code. This difference is obtained by counting the number of direction changes (in, say, counter clockwise direction) that separate two adjacent element of code. Size normalization can be achieved by altering the size of the re-sampling grid.

### 3.5 Feature Extraction using Density, Moment, Fourier Descriptor and Chain Code

1. **Feature Set I:** Density features based upon the zoning approach as described in section 3.1 is calculated. The result is stored as feature vector of size is 65.
2. **Feature Set II:** Augmenting the density features (as described in section 3.1) with central moment features (as described in section 3.2) of order three and four for the numeral image to get feature vector of size 161.
3. **Feature Set III:** Feature vector comprising 64 dimensional Fourier descriptors invariant to rotation, scale, and translation. Procedure is described in section 3.2.

4. **Feature Set IV:** Feature vector comprising chain code and Fourier Descriptors of magnitude 108 is used as feature. Procedure is describe in section 3.3
5. **Feature Set V:** Features computed by normalizing the chain code and augmenting it with Fourier features. The size of the feature vector is 48. Procedure is described in section 3.3 and 3.4.

## 4. Classification

We have used SVM classifier in present system. The SVM classifier is basically a two class classifier based on the discriminate functions. A discriminate function represents a surface, which separates the patterns as two classes. For OCR applications a number of two class classifiers are trained with each one distinguishing one class from the other. Each class label has an associated SVM and a test example is assigned to the label of the class whose SVM gives the largest positive output.

## 5. Experimental Result

The effectiveness of the features proposed in the paper are evaluated by performing experiments on our own database containing 12690 isolated Marathi handwritten numeral images obtained from writers belonging to different age groups and professions. Feature set V comprising 32 dimensional Fourier descriptors and normalized chain code of length 16, has yielded highest recognition accuracy of 98.15 %. The Table 1 provides comparative average recognition rates obtained by using different feature sets.

Table 1: Recognition accuracy using SVM

Input Num.	Density Features using zoning	Density and Central Moments	FD	FD and Chain Code	FD and Normalized chain Code
0	100.00	100.00	99.61	99.69	99.92
9	96.15	97.16	96.77	97.08	98.74
2	95.60	96.22	97.16	97.01	97.01
3	96.78	97.72	97.09	97.56	98.03

8	98.27	98.19	97.16	97.08	97.79
5	96.07	96.07	98.35	96.54	98.03
3	96.46	97.72	96.69	97.32	96.69
9	98.90	98.97	96.21	98.26	97.71
7	99.29	99.37	98.66	98.74	98.90
6	97.01	97.48	97.71	98.58	98.66
Avg.	97.45	97.89	97.54	97.79	98.15

The error rate for five different proposed methods is shown in Figure 2. The error rate for numerals 0, 4, 8, and 9 is very less compared to error rates for other numerals. The methods proposed in this thesis were implemented using Mat-lab 7.0 software and PR Tools on Pentium systems.

## 6. Conclusion

The performance of the method based upon 'Fourier Descriptor and Normalized Chain Code', presented in this paper better as compared with other methods. The proposed method performs well and appears promising as compared to other methods in the literature.

**Table 2 Performance comparisons of proposed method with other methods in literature**

Method	Data	Features	Feat- ure	Class- ifier	Reco- gnition
M. Hanma- ndlu et.al. [35]	3500	Normalize distance	48	Fuzzy	96.00 %
R. J. Ramteke et.al. [47]	2000	Invariant moments	78	Gaussian	92.28 %
Reena bajaj [52]	2460	Density, Moment	48	MLP	89.68 %

U. Pal et.al. [74]	22546	Chain Code	64	Quadratic	98.86 %
P. M. Patil et.al. [45]	2000	Ring data	---	Fuzzy	99.5 %
Benne R.G. et.al. [8]	1500	Water Reservoir etc.	13	k-NN	90.20 %
Proposed	12690	FD and Chain code	48	SVM	98.15 %

## 7. References

- [ 1 ] Mowlaei and K. Faez, "Recognition of Isolated Handwritten Persian/Arabic Characters and Numerals Using Support Vector Machines", Proceedings of XIII Workshop on Neural Networks for Signal Processing, pp. 547-554, 2003.
- [ 2 ] Alireza Alaei, Umapada Pal and P. Nagabhushan, "Using Modified Contour Features and SVM Based Classifier for the Recognition of Persian/Arabic Handwritten Numerals", Proceedings of International Conference on Advances in Pattern Recognition, pp.391-394, 2009.
- [ 3 ] Banashree N. P., and R. Vasanta, OCR for Script Identification of Hindi (Devnagari) Numerals using Feature Sub Selection by Means of End-Point with Neuro-Memetic Model, International Journal of Intelligent Systems and Technologies 2;3, pp 206-210, 2008.
- [ 4 ] Benne R.G., Dhandra B.V.and Mallikarjun Hangarge, "Tri-scripts handwritten numeral recognition: a novel approach", Advances in Computational Research, Volume 1, Issue 2, pp-47-51, 2009.
- [ 5 ] F. Kimura, T. Wakabayashi, S.Tsuruoka and Y. Miyake, "Improvement of handwritten Japanese character recognition using weighted direction code histogram", Pattern recognition, vol.30, no.8, pp.1329-1337, 1997.
- [ 6 ] G. S. Lehal and Nivedan Bhatt, "A Recognition System For Devnagri And English Handwritten Numerals ", Advances in Multimodal Interfaces – ICMI 2001, T. Tan, Y. Shi and W. Gao (Editors), Lecture Notes in Computer Science, Vol. 1948, Springer-Verlag,Germany, pp. 442-449. (2000).
- [ 7 ] G. G. Rajaput and Mallikarjun Hangarge, "Recognition of isolated handwritten Kannada numerals based on image fusion method: ",PREMI07, LNCS.4815, pp.153-160, 2007.

- [ 8 ] H. Khosravi, E. Kabir, "Introducing a very large dataset of handwritten Farsi digits and a study on the variety of handwriting styles", *Pattern Recognition Letters* Vol.28, Issue 10, pp. 1133-1141, 2007.
- [ 9 ] H. Soltanzadeh and M. Rahmati, "Recognition of Persian handwritten digits using image profiles of multiple orientations", *Pattern Recognition Letters* 25, pp.1569-1576, 2004.
- [ 10 ] Javad Sadri Ching Y. Suen Tien D. Bui, "Application of Support Vector Machines for Recognition of Handwritten Arabic/Persian Digits", 2ndMVIP, K.N. Toosi Univ. of Tech., Tehran, Iran. Feb. 2003. Vol. 1, pp.300-307, Feb. 2003.
- [ 11 ] M. H. Shirali-Shahreza, K. Faez and A. Khotanzad, "Recognition of Hand-written Persian/Arabic Numerals by Shadow Coding and an Edited Probabilistic Neural Network", *Proceedings of International Conference on Image Processing*, Vol. 3, pp. 436-439, 1995.
- [ 12 ] M. Hanmandlu, A.V. Nath, A.C. Mishra and V.K. Madasu, "Fuzzy Model Based Recognition of Handwritten Hindi Numerals using Bacterial Foraging", 6th IEEE/ACIS International Conference on Computer and Information Science (ICIS 2007), Computer Society, 2007.
- [ 13 ] M. Hanmandlu, J. Grover, V. K. Madasu and S. Vasikarla, "Input fuzzy for the recognition of handwritten Hindi numeral.", *International Conference on Informational Technology*, vol. 2, pp. 208-213., 2007.
- [ 14 ] M. Ziaratban, K. Faez and F. Faradji, "Language-Based Feature Extraction Using Template-Matching in Farsi/Arabic Handwritten Numeral Recognition", *Proceedings of 9th ICDAR*, Vol.1, pp. 297-301, 2007.
- [ 15 ] P.M. Patil, T.R. Sontakke, "Rotation scale and translation invariant handwritten Devanagiri numeral character recognition using fuzzy neural network", Elsevier, *Pattern Recognition*, vol. 40, pp. 2110-2117. 2007.
- [ 16 ] R. J. Ramteke and S. C. Mehrotra, "Recognition of Handwritten Devnagari Numerals", *International Journal of Computer Processing of Oriental Languages Ó Chinese Language Computer Society & World Scientific Publishing Company*, 2008.
- [ 17 ] Reena Bajaj, Lipika Dey and Santanu Chaudhur, "Devnagari numeral recognition by combining decision of multiple connectionist classifiers", Vol. 27, Part 1, pp. 59-72, February 2002.
- [ 18 ] S.V. Rajashekararadhya and Dr P. Vanaja Ranjan, "Efficient zone based feature extration algorithm for handwritten numeral recognition of four popular south indian scripts", *Journal of Theoretical and Applied Information Technology*, pp. 1171-1180, 2005-06.
- [ 19 ] S.V. Rajashekararadhya and P. Vanaja Ranjan, "Efficient handwritten numeral recognition of Kannada and Telgu scripts", *International conference on Sensors, Security, Software and Intelligent System (ISSIS)*, pp.19-23, 2009.
- [ 20 ] S.V. Rajashekararadhya and P. Vanaja Ranjan, "Zone based Feature Extraction Algorithm for Handwritten Numeral Recognition of Kannada Script, Advance Computing Conference, 2009, IEEE Intl., pp 525-528.
- [ 21 ] U. Pal, N. Sharma, T. F. Kimura, S. Pal, "Recognition of Off-Line Devnagari Characters Using Quadratic Classifier", *Lecture Notes in Computer Science*, 2006, Volume 4338, pp. 805-816, 2006.