# Gujarati handwritten numeral optical character reorganization through neural network

Apurva A. Desai *

*Veer Narmad South Gujarat University, Surat, Gujarat, India*

## ARTICLE INFO

## ABSTRACT

This paper deals with an optical character recognition (OCR) system for handwritten Gujarati numbers. One may find so much of work for Indian languages like Hindi, Kannada, Tamil, Bangala, Malayalam, Gurumukhi etc, but Gujarati is a language for which hardly any work is traceable especially for handwritten characters. Here in this work a neural network is proposed for Gujarati handwritten digits identification. A multi layered feed forward neural network is suggested for classification of digits. The features of Gujarati digits are abstracted by four different profiles of digits. Thinning and skew-correction are also done for preprocessing of handwritten numerals before their classification. This work has achieved approximately 82% of success rate for Gujarati handwritten digit identification.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

Gujarati belonging to Devnagari family of languages, which originated and flourished in Gujarat—a western state of India, is spoken by over 50 million people of the state. Though it has inherited rich cultural and literature, and is a very widely spoken language, hardly any significant work has been done for the identification of Gujarati optical characters. The Gujarati script differs from those of many other Indian languages not having any shirolekha (headlines). Gujarati numerals do not carry shirolekha and it applies to almost all Indian languages. The numerals in Indian languages are based on sharp curves and hardly any straight lines are used. Fig.1 is a set of Gujarati numerals.

Moreover, some of the digits in various languages share relatively similar shapes although they have different meanings in different languages. Table 1 presents the digits 0 – 9 as they appear in many other Indian languages.

As it is visible in Fig. 1, Gujarati digits are very peculiar by nature. Only two Gujarati digits one (1) and five (5) are having straight line, making Gujarati digit identification a little more difficult. Also Gujarati digits often invite misclassification. These confusing sets of digits are as shown in Fig. 2.

As shown in Fig. 2, digits zero (0, three (3) and seven (7), digits one (1) and six (6), and digits eight (8) and nine (9) share similar shapes.

This paper addresses the problem of handwritten Gujarati numeral recognition. Gujarati numeral recognition requires smoothing, thinning, skew detection and correction, and normalization process is performed as preprocess. Further, profiles are used for feature extraction and artificial neural network (ANN) is suggested for the classification. This paper is organized in the following sections; Section 2 describes the early attempts in Indian language OCR. Section 3 explains suggested preprocessing. Section 4 is devoted to feature extraction. Section 5 describes the suggested artificial neural network and Section 6 describes the training of the network. Lastly in Section 7 the results are explained. The result section is then followed by conclusions.

## 2. Some early attempts in Indian language OCR

A lot of research work has been done and is still being done in OCR for various languages. More and more researchers are attracted to this challenging field. Though this technology has advanced to a much high level, for English that is not the case for many Indian languages. Since not much work has been done on Gujarati and as some similarities are found in the patterns of numbers of various Indian languages as evident in Table 1, I have taken the research work done in other Indian languages as a reference point. Further, since OCR is categorized in two classes,

* Correspondence address: Department of Computer Science, Veer Narmad South Gujarat University, Udhna Magdalla Road, Surat 395 007, Gujarat, India. Tel.: +91 261 2257906.
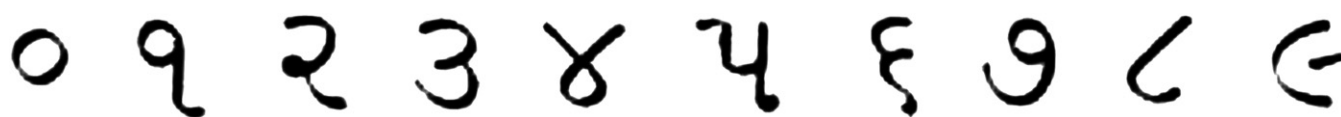*E-mail address:* desai_apu@hotmail.com
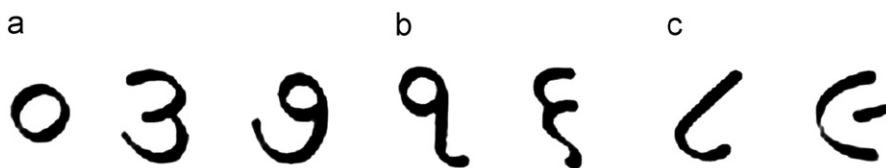
**Fig. 1.** Gujarati digits 0–9.



**Fig. 2.** Confusing Gujarati digits.

**Table 1**
Digits 0–9 in various Indian languages.

| Variant | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| Gujarati | ૦ | ૧ | ૨ | ૩ | ૪ | ૫ | ૬ | ૭ | ૮ | ૯ |
| Gurumukhi | ੦ | ੧ | ੨ | ੩ | ੪ | ੫ | ੬ | ੭ | ੮ | ੯ |
| Kannada | ೦ | ೧ | ೨ | ೩ | ೪ | ೫ | ೬ | ೭ | ೮ | ೯ |
| Malayalam | ൦ | ൧ | ൨ | ൩ | ൪ | ൫ | ൬ | ൭ | ൮ | ൯ |
| Telugu | ౦ | ౧ | ౨ | ౩ | ౪ | ౫ | ౬ | ౭ | ౮ | ౯ |

OCR for printed characters and OCR for handwritten characters, the literature is discussed accordingly.

One can trace the history of printed character recognition in Indian languages since 1977. I could find the first ever, traceable attempt made by Rajsekaran et al. [1] in the year. They presented their work to identify printed Telugu characters, and it identified 50 printed primitive Telugu characters. This work used the shapes of characters for identification, followed by the statistical approach, the decision tree, for classification. After this paper one will find various contributory works in Indian languages by many researchers. Sinha et al. [2] in 1979 presented work for Devnagari scripts. In this work they analyzed the structural characteristics of Devnagari scripts. Here various aspects of Devnagri scripts were also reported. In 1984 Ray et al. [3] presented work on Bengala character recognition based on nearest neighbor classifier. Rao et al. [4] used feature-based approach to recognize Telugu scripts in 1995, suggesting the use of canonical shapes of printed Telugu characters, and achieving success rate between 78% and 90%. The use of neural network can be found for recognition of characters of an Indian language in work presented by Sukhswami et al. [5] in 1995. In this work they recognized Telugu characters. Also authors suggested multiple neural network associative memory (MNNAM) model to handle

the problem of shortage of memory. In the year 1998 Chaudhari et al. [6] presented a work to recognize printed Bangala characters, using templates matching and feature-based tree classifier. The authors also suggested characters unigram statistics to make the tree classifier more efficient. In this work 95.5% of success was recorded. The first ever work of character recognition for printed or digitized Gujarati language was done in 1999. In this year Antani et al. [7] presented a work in which they used euclidean minimum distance classifier (EMDC) and hamming distance classifier (HDC) to classify various printed Gujarati characters. However they obtained a very low recognition rate of 67%. In 2002, Garain et al. [8] proposed an algorithm, based on fuzzy multi-factorial analysis for segmentation of touching printed Devnagari and Bangala scripts. Chaudhuri et al. [9] presented a work on the recognition of printed Oriya script in 2002. This work presented line segmentation, word segmentation, character recognition and character segmentation. It achieved on average 97% accuracy. Chakravarthi et al. [10] noted that it was difficult to achieve more than 93% of recognition accuracy. They have also listed out some of the factors to improve recognition accuracy up to 97%. Lakshmi et al. [11] presented their work to recognize printed basic symbols of Telugu language. This work, presented in 2003, used seven moments for feature abstraction

and then K-nearest neighbor algorithm. In this work authors have reported 90–100% accuracy. Many times in a single document more than one script are used. This situation was considered by Pal et al. [12] in 2003. In this work they identified different lines printed in different scripts. They did this with the help of characteristics of various scripts, contour tracing, water reservoir principle, left and right profiles etc. For classification of languages they have made use of binary tree classifier. Pal et al. [13] presented a survey of character recognition in 2004. Dholakia et al. [14], in 2005, presented their work on Gujarati language. In this work they presented an algorithm to identify various zones used for Gujarati printed text. Here authors have used horizontal and vertical profiles for line and word segmentation. Further, the zones are identified by slope of lines created by upper left corner of rectangle created by the boundaries of connected components. In this work 95% of accuracy has been recorded for segmentation. In 2006, Anuradha et al. [15] presented their work for digitization of old documents. When Indian scripts are used, horizontal segmentation becomes very difficult. Jindal et al. [16] addressed this problem for eight major Indian printed scripts in 2007. They achieved 96.45–99.79% of accuracy through this presented algorithm. In 2008, Anuradhasrinivas et al. [17] gave an exhaustive overview of OCR research in Indian scripts. In this paper they have surveyed in detail various Indian languages. In the same year Manjunatharadhya et al. [18] presented a system for multilingual character recognition for printed south Indian scripts, Kannada, Telugu, Malayalam and Tamil. Here they have used Fourier transform and principal component analysis (PCA). Good efficiency has been recorded in this work for both printed document and also tested for handwritten characters.

Compared to OCR for printed characters, very limited work can be traced for handwritten character recognition for Indian languages. Chinnuswami et al. [19] in 1980 presented their work to recognize hand printed Tamil Characters. In this paper, authors have used structure of the Tamil characters. Using the curves and strokes of characters, the features were identified and statistical approach was used for the classification. Dutta et al. [20] recognized both printed and handwritten alpha–numeric Bengali characters using curvature features in 1993. Here features like curvature maxima, curvature minima, and inflexion points were considered. In this work recognition was performed on isolated characters. Thinning and smoothing were also performed prior to classification of characters. In 2007, Banashree et al. [21] attempted identification of handwritten Hindi digits, using diffusion half toning algorithm. 16-segment display concept has been used here for feature extraction. They proposed a neural classifier for classification of isolated digits. Here they achieved accuracy level up to 98%. In 2008, Rajashekararadhya [22] proposed an offline handwritten OCR technique for four south Indian languages like Kannada, Telugu, Tamil and Malayalam. In this work they suggested a feature extraction technique, based on zone and image centroid. They used two different classifiers nearest neighborhood and back propagation neural network to achieve 99% accuracy for Kannada and Telugu, 96% for Tamil and 95% for Malayalam. In 2009, Shanthi et al. [23] used support vector machine (SVM) for handwritten Tamil characters. In this work image subdivision was used for feature extraction. They recorded 82.04% accuracy.

## 3. Preprocessing

For developing a system to identify Gujarati handwritten digits, I have collected numerals 0–9 written in Gujarati scripts from 300 different people of various background and different genders. These numbers were scanned in 300 dpi by a flatbed scanner. Initially they are in separate boxes of $90 \times 90$ pixels each. Since our problem is to identify handwritten digits, the first thing required is to bring all the characters in a standard normal form. This is needed because when a writer writes he may use different types of pens, papers, they may follow even different styles of writing etc. In Fig. 3 the scanned images are illustrated which show the need for preprocessing of the images before processing them for the classification.

In Fig. 3a the set of characters shows digits written on a blotting paper type of paper and that too with an ink pen. It is notable about in these digits that the boundaries of these characters are very uneven and smoothness is absent there. Also because of the spreaded ink the thickness of characters in Fig. 3a is more than that in Fig. 3b. Fig. 3b is a set of scanned numerals which are written on a good quality paper. You can see the digits in Fig. 3b are thinner and having smoother boundaries. Fig. 3c is an interesting case. Most of the people keep their paper in an angular position while writing. As a result a skewness is created in most of the characters. Observe this skewness in digit five which is shown in Fig. 3c and compare it with the first digit of Fig. 3c which is a printed digit and that is why there is no skewness in it. Next characteristic which needs to be addressed is the size. When handwritten digits appear in picture we need to deal with the digits with different sizes as all the people write in different sizes. Fig. 3d shows a single digit, six, written in different sizes by different people.

The first preprocess which is performed is an adjustment of the contrast. The contrast adjustment is required to remove the use of ink of different colors and also different intensities of black. To adjust contrast the contrast limited adaptive histogram equalization (CLAHE) algorithm is employed here in this work. This algorithm is utilized with $8 \times 8$ tiles and 0.01 contrast enhancement constant with uniform distribution. Further, the boundaries of images of digits are required to be smoothening out. The smoothing is done using median filter. The median filter is performed on each $3 \times 3$ neighborhood pixels. Fig. 4 shows the digits before and after contrast adjustment and smoothing process.

Further the digits are in different sizes. For OCR, we need to put all the handwritten digits in an uniform size that is—the digits should be in normal form. Ashwin et al. [24] have proposed the use of support vector machine (SVM) method for size independent identification of Kannada printed characters. In order to keep our algorithm simple, here in this work all the digits are reconstructed in the size of $16 \times 16$ pixels, using nearest neighborhood interpolation (NNI) algorithm.
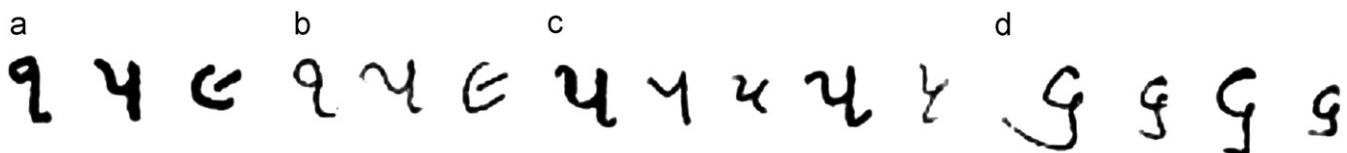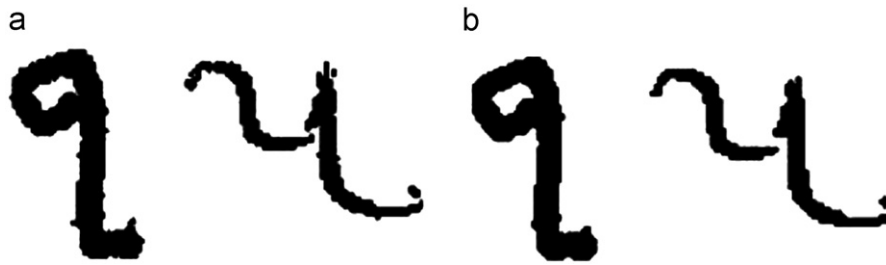


**Fig. 3.** Scanned digits of different nature.

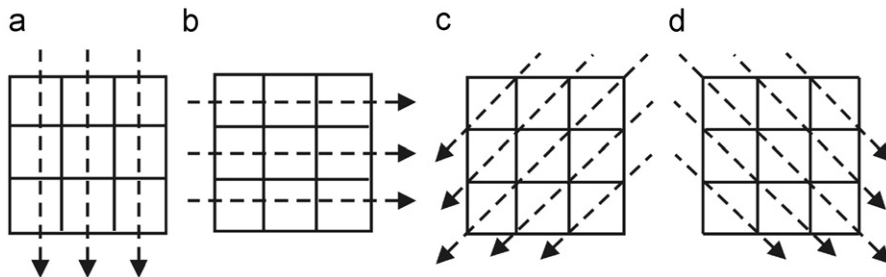**Fig. 4.** Smoothing process: (a) original images; (b) smooth images.



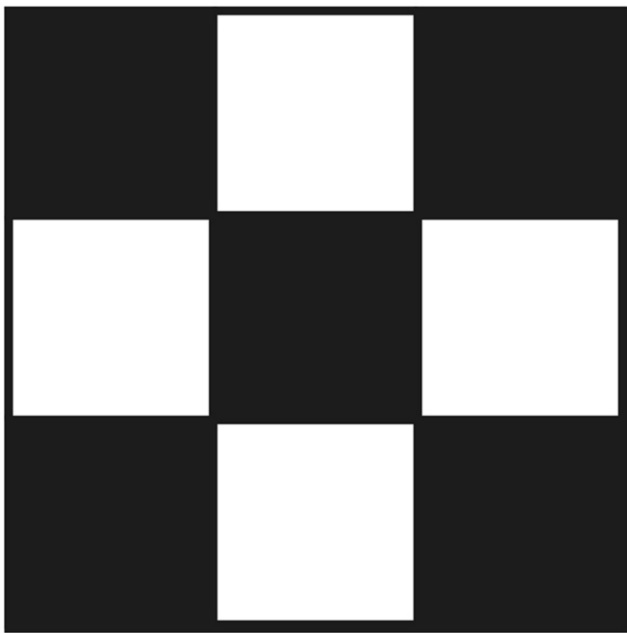**Fig. 5.** Pattern profile of $3 \times 3$ pattern matrix.



**Fig. 6.** $3 \times 3$ Pattern.

## 4. Feature extraction

In optical character recognition, the most important aspect is extraction of features of each of the numerals that are required to be identified. As seen in Section 1, most of the Gujarati numerals are based on very sharp curves, and for handwritten numerals also these curves may be very irregular. In this case it is very difficult to extract different objects or components which constitute unique features for each of the individual Gujarati numerals. To handle this situation here in this work, various profiles of digits are used as template to identify various digits. In this very simple but effective, feature extraction technique the use of four different profiles, horizontal, vertical, and two diagonals, is

suggested. It is like a puzzle to solve. For example, if we take a $3 \times 3$ box and you are asked to activate all pixels which has 2,1 and 2 activated pixels in the first, second and third row, respectively, this is the X-profile of the pattern. In Fig. 5 such four profiles are shown for a $3 \times 3$ pattern.

Thus if the four profiles—X-profile, Y-profile, diagonal1 profile and diagonal2 profile—for a pattern is (2,1,2), (2,1,2), (1,0,3,0,1) and (1,0,3,0,1) respectively, the pattern in question would be as shown in Fig. 6.

The vector of these four profiles, that is [X-profile, Y-profile, diagonal1 profile, diagonal2 profile], can be used as an abstracted feature for identification of a digit. For the pattern shown in Fig. 6, the feature vector or profile vector is [2,1,2,2,1,2,1,0,3,0,1,1,0,3, 0,1]. For Gujarati numeral recognition profile vector is created for all the digits which are converted into $16 \times 16$ pixels after preprocessing. The orientation or the skew which exists in the writing of most of the writers is the next thing to remove from the digits before their identification. Chaudhuri et al. [25] have proposed an algorithm to remove skew from digitized Indian script. But in this work, skew is detected and removed from the entire document. Here we need to deal with skew correction for individual digits. To handle this different situation of skew, rotation up to $10°$ about center point $(x_c, y_c)$ of the digit is performed on each of the standard digit which are considered in the train set. Thus five more patterns for each of the digits are created in both clock wise and anti clock wise directions with the difference of $2°$ each. The following geometric transformations (1) and (2) are performed on each and every pixel for the clock wise and anti clock wise direction, respectively,

$$x' = x\cos\theta + y\sin\theta + x_c(1-\cos\theta) - y_c\sin\theta$$

$$y' = -x\sin\theta + y\cos\theta + y_c(1-\cos\theta) + y_c\sin\theta \qquad (1)$$

$$x' = x\cos\theta - y\sin\theta + x_c(1-\cos\theta) + y_c\sin\theta$$

$$y' = x\sin\theta + y\cos\theta + y_c(1-\cos\theta) - x_c\sin\theta \qquad (2)$$

Here, $(x', y')$ is the pixel after transformation, $(x, y)$ is pixel before transformation and the $\theta$ is the angle of rotation. Fig. 7
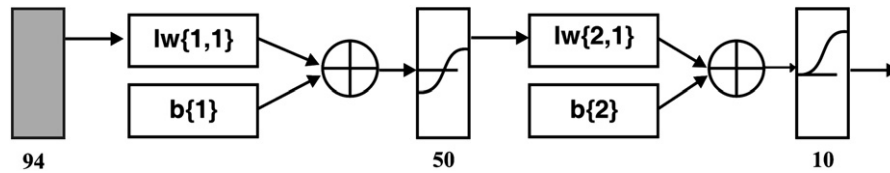
**Fig. 7.** Digit 5 in different orientation.



**Fig. 8.** Feed forward back propagation neural network.

**Table 2**
Neural network details.

| Network | Feed forward back propagation | |
|---|---|---|
| Layers | 3 (Input, Hidden, Output) | |
| | | |
| Layer details | Neurons | Training function |
| Input layer | 94 | transig |
| Hidden layer | 50 | transig |
| Output layer | 10 | logsig |
| | | |
| Back propagation training function | traingdx | |
| Back propagation bias learning function | learngdm | |
| Performance function | MSE | |
| Performance goal | 0.01 | |
| Input range | [0.94] | |

shows a resultant images of a digit which can then be considered in the training set of a digit 5 after rotation.

## 5. The neural network

As Luh Tan et al. [26], Sukhswami et al. [27], Wellner et al. [28] etc. have used neural network for character classification, in this work of Gujarati numeral classification too, neural network is suggested. A feed forward back propagation neural network is used for Gujarati numeral classification. This proposed multi-layered neural network consists of three layers with 94, 50, and 10 neurons, respectively. The input layer is the layer which accepts the profile vector which is of $1 \times 94$ in size. As this network is used for classification of 10 digits, it has 10 neurons in the output layer. Fig. 8 shows this feed forward back propagation neural network.

In Table 2 further specification details of the network used for classification of Gujarati numbers are given.

## 6. The training of network

For this experiment, a total of 278 responses were taken into consideration. Thus total images considered for this experiment is 2798, taking 10 digits per responder. Out of these 278 sets of various images of digits, 11 sets were created by a standard font. For training, the features are abstracted first for all of these images of digits. A profile vector for a digit five is shown here

[1 2 1 1 1 7 8 1 1 1 1 2 15 10 3 2 6 5 4 4 4 4 4 5 7 1
1 1 1 4 4 2 0 0 0 0 0 0 0 0 0 0 0 0 4 6 8 7 6 5 5
2 2 2 2 2 2 2 2 0 0 0 0 2 3 2 2 1 1 1 1 1 2 3 4 3 3
3 3 2 3 2 2 2 2 2 1 0 1 1 2 1 1]

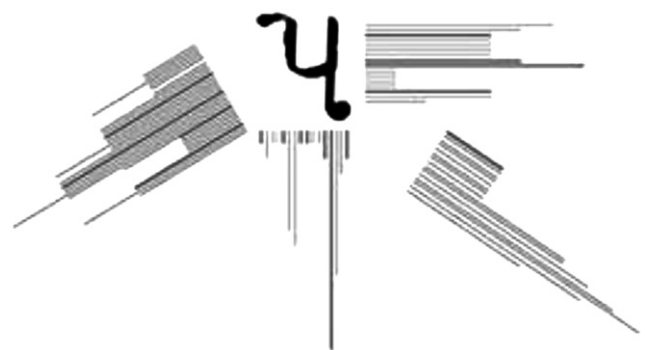Fig. 9 is demonstrating how this profile vector represents the feature of the same digit.



**Fig. 9.** Features represented by a profile vector.

Initially, the neural network is trained for one set of standard digits zero to nine and 10 sets of skew corrected digits zero to nine; i.e. the network is initially trained for 110 images which were created using a standard font. After the training this network is further trained for more 50 sets of handwritten digits that are for 500 digits. Fig. 10 shows the performance graph of one of such training sessions.

In the Fig. 11 the complete process of Gujarati numeral optical character recognition is shown.

## 7. Results

As mentioned above this network was trained for total 61 sets of digits, and was tested for 203 other new sets of digits and also for the sets by which the network was trained. In total the network was trained by 610 digits and tested for 2650 digits.

In totality this network gave 81.66% of success rate. The success rate for the training set is higher than the overall success rate. The success rate of the first 110 digits is just 71.82% but for the later 500 digits is 91.0 %. The results are summarized in Table 3.

Let us examine the results obtained for the different digits. The Table 4 shows the confusion among, the identified digits while testing the proposed network for the training sets of handwritten digits.

While testing the network, it is seen that the success rate of identification of zero, four and seven, is very high, i.e. 98%. And it is as low as 72.0% for the digit six. The digit six is misidentified as three for seven times out of fifty, that means six is identified as three for 14% of total attempts. Similarly the digit one is identified as four for 6% attempts. The digit three is misidentified as seven for 14% of attempts. Thus the most confusing digit here is the digit 6. Table 5 shows success rate for all 2650 handwritten digits. This table shows more interesting results. The maximum success rate achieved is for the digit four and that is 96.23%. The success rate is also high for digits eight and zero and that is 92.46% and 91.70%, respectively. The success rate of this proposed network is very low for the digits six and three with 60.0% and 67.55%, respectively. Digit zero has been misread as nine and seven for 3.78% and 3.02%, respectively. Zero is also misidentified as one, two, four and eight for one attempt each. Digit two is misidentified as digit four for 32 attempts out of 265, i.e. for 12.08%. Digit one is also misidentified as digit five for 3.02%. One is also misidentified as two, six and three. Digit one has never been identified as three and seven, Digit two is misidentified as
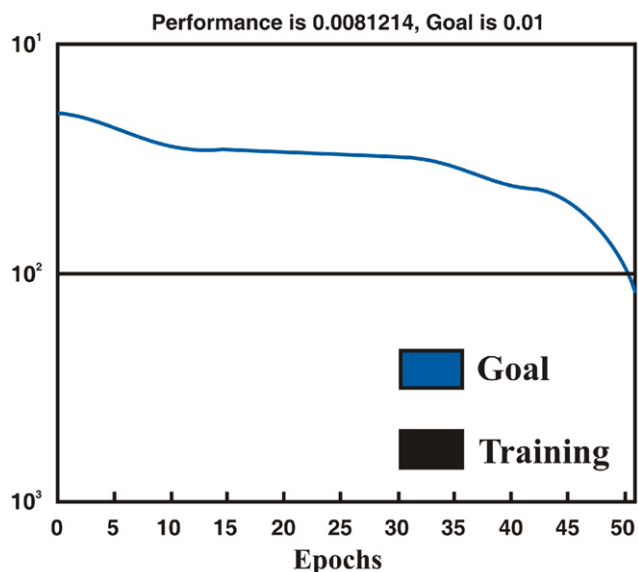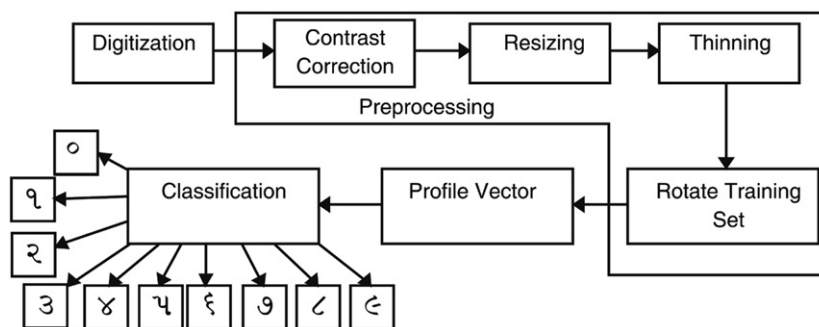


Fig. 10. Performance of training.



Fig. 11. Gujarati handwritten optical character recognition process.

**Table 3**
Result summary.

| Sets | Nos of digits | Type of sets | Success rate (%) |
|---|---|---|---|
| 11 sets of standard fonts | 110 | Training sets | 71.82 |
| 50 sets of handwritten digits | 500 | Training sets | 91.0 |
| 200 sets of handwritten digits | 2000 | Test sets | 81.50 |
| Total 265 sets | 2650 | | 81.66 |

**Table 4**
Performance of network for test set.

| Digits | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Success (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 49 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 98.0 |
| 1 | 0 | 44 | 1 | 0 | 3 | 1 | 1 | 0 | 0 | 0 | 88.0 |
| 2 | 0 | 1 | 48 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 96.0 |
| 3 | 0 | 0 | 0 | 40 | 3 | 0 | 0 | 7 | 0 | 0 | 80.0 |
| 4 | 0 | 1 | 0 | 0 | 49 | 0 | 0 | 0 | 0 | 0 | 98.0 |
| 5 | 0 | 1 | 0 | 0 | 2 | 46 | 0 | 0 | 1 | 0 | 92.0 |
| 6 | 4 | 1 | 0 | 7 | 0 | 0 | 36 | 2 | 0 | 0 | 72.0 |
| 7 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 49 | 0 | 0 | 98.0 |
| 8 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 48 | 0 | 96.0 |
| 9 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 3 | 46 | 92.0 |

**Table 5**
Performance of network for all the digits.

| Digits | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Success (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 243 | 1 | 1 | 0 | 1 | 0 | 0 | 8 | 1 | 10 | 91.70 |
| 1 | 1 | 212 | 5 | 0 | 32 | 8 | 3 | 0 | 3 | 1 | 80.00 |
| 2 | 0 | 4 | 207 | 0 | 13 | 25 | 8 | 1 | 4 | 3 | 78.12 |
| 3 | 4 | 0 | 1 | 179 | 16 | 1 | 8 | 55 | 1 | 0 | 67.55 |
| 4 | 0 | 5 | 0 | 1 | 255 | 2 | 0 | 0 | 2 | 0 | 96.23 |
| 5 | 0 | 16 | 0 | 0 | 30 | 204 | 4 | 6 | 4 | 1 | 76.99 |
| 6 | 22 | 12 | 7 | 30 | 16 | 5 | 159 | 13 | 0 | 1 | 60.00 |
| 7 | 10 | 3 | 0 | 7 | 6 | 1 | 1 | 235 | 1 | 1 | 88.68 |
| 8 | 2 | 6 | 3 | 0 | 1 | 1 | 0 | 0 | 245 | 7 | 92.46 |
| 9 | 5 | 2 | 10 | 0 | 1 | 0 | 0 | 0 | 22 | 225 | 84.91 |

five for 9.44% and as four for 4.91%. It is misidentified as six for eight times and four times each for digits one and eight. Digit two is never identified as digit zero and three. The failure rate of identification of digit three is maximal for digit seven and that is 20.76% and for digit four it is also as high as 6.04%. Digit four is confused as digit one in five attempts out of total 265 tests. Digit five is again a confusing digit for this proposed network. In 11.32% five is misidentified as four and 6.04% as one. Digit five is never identified as zero and two. Digit six has given weakest performance. It is misidentified as three for 11.32%, as zero for 8.31% and as four for 6.04%. Six is never identified as eight. Digit seven is identified as zero in ten attempts and in seven attempts is identified as three. Network has misidentified digit eight as nine, one and two in 7, 6 and 3 attempts, respectively, out of 265. The failure rate of digit nine is high for digit eight and it is 8.31%.

In Section 2, we have noted down the accuracy level that has been obtained by various researchers for handwritten character recognition in some of the Indian languages. The result that has been achieved here in this work for handwritten Gujarati numeral recognition, i.e. about 82% of success rate, is very good compared to the even printed character recognition [7] in Gujarati. There is no other literature that I could find on handwritten Gujarati optical characters recognition for comparison. Hindi numeral [21] has demonstrated good accuracy which is higher than the work presented here for Gujarati numerals. Compared to handwritten character recognition for the Kannada, Telugu, Tamil, Malayalam languages [22,23], the result obtained here is low. However, the result obtained for Tamil character recognition [23] is as good as results presented in this work.

## 8. Conclusion

In this work a feed forward back propagation neural network is proposed for the classification of the Gujarati numerals. Various techniques are used in the preprocessing phase before implementing classification of numerals. The overall performance of this proposed network is as high as 81.66%, but still it is not up to the mark. The performance of any classification model is mainly based on the feature abstraction. To improve the performance of this prototype, the improved features abstraction technique and/or the preprocessing techniques are possibly required. As a whole, this model offers a satisfactory success rate but it is subject to further improvement.

## References

[1] S.N.S. Rajasekaran, B.L. Deekshatulu, Recognition of printed Telugu characters, Computer Graphics Image Processing (1977) 335–360.

[2] R.K. Sinha, H.N. Mahabala, Machine recognition of Devnagari script, IEEE Transactions on Systems, Man, and Cybernetics (1979) 435–441.

[3] K. Ray, B. Chatterjee, Design of a nearest neighbor classifier system for Bengali character recognition, Journal of Institute of Electronics, Telecom Engineers 30 (1984) 226–229.

[4] P.V.S. Rao, T.M. Ajitha, Telugu script recognition—a feature based approach, Proceedings of ICDAR, IEEE (1995) 323–326.

[5] R. Sukhaswami, P. Seetharamulu, A.K. Pujari, Recognition of Telugu characters using neural networks, International Journal of Neural System 6 (1995) 317–357.

[6] B.B. Chaudhuri, U. Pal, A complete printed Bangla OCR system, Pattern Recognition 31 (1998) 531–549.

[7] S. Antani, L. Agnihotri, Gujarati character recognition, in: Fifth International Conference on Document Analysis and Recognition (ICDAR'99), 1999, pp. 418–421.

[8] U. Garain, B.B. Chaudhuri, Segmentation of touching characters in printed Devnagari and Bangla Scripts using fuzzy multifactorial analysis, IEEE Transactions on Systems, Man and Cybernetics, Part C 32 (4) (2002) 449–459.

[9] B.B. Chaudhhuri, U. Pal, M. Mitra, Automatic recognition of printed Oriya script, Saadhanaa 27 (1) (2002) 23–34.

[10] B. Chakravarthy, T. Ravi, S.M. Kumar, A. Negi. On developing high accuracy OCR systems for Telugu and other Indian scripts, in: Proceedings of Language Engineering Conference, 2002, pp. 18–23.

[11] C.V. Lakshmi, C. Patvardhan, Optical Character recognition of basic symbols in printed Telugu text, IE(I)Journal-CP 84 (2003) 66–71.

[12] U. Pal, S. Sinha, B.B. Chaudhuri, Multi script line identification from Indian documents, in: Proceedings of the Seventh International Conference on Document Analysis and Recognition (ICDAR 2003), 2003, pp. 880–884.

[13] U. Pal, B.B. Chaudhuri, Indian script character recognition: a survey, Pattern Recognition 37 (2004) 1887–1899.

[14] J. Dholakia, A. Negi, S. Ram Mohan, Zone identification in the printed Gujarati text, in: Proceedings of the 2005 Eight International Conference on Document Analysis and Recognition (ICDAR05), 2005.

[15] B. Anuradha, B. Koteswarrao, An efficient Binarization technique for old documents, in: Proceedings of International Conference on Systemics, Cybernetics, and Inforrmatics (ICSCI2006), Hyderabad, 2006, pp. 771–775.

[16] M.K. Jindal, R.K. Sharma, G.S. Lehal, Segmentation of horizontally overlapping lines in printed Indian scripts, International Journal of Computational Intelligence Research 3 (4) (2007) 277–286.

[17] B. Anuradhasrinivas, A. Agarwal, R. Rao, An overview of OCR research in indian scripts, International Journal of Computer Science and Engineering System 2 (2008) 141–153.

[18] V.N. Manjunatharadhya, G. Hemanthakumar, S. Naushath, Multilingual OCR system for South Indian scripts and English documents: an approach based on Fourier transform and principal component analysis, Engineering Application of Artificial Intelligence 21 (4) (2008) 658–668.

[19] P. Chinnuswamy, S.G. Krishnamoorty, Recognition of hand-printed Tamil characters, Pattern Recognition 12 (3) (1980) 141–152.

[20] A. Dutta, S. Chaudhary, Bengali Alpha-numeric character recognition using curvature features, Pattern Recognition 26 (12) (1993) 1757–1770.

[21] N.P. Banashree, D. Andhre, R. Vasanta, P.S. Satyanarayana, OCR for script identification of Hindi (Devnagari) numerals using error diffusion Halftoning Algorithm with neural classifier, Proceedings of World Academy of Science Engineering and Technology 20 (2007) 46–50.

[22] S.V. Rajashekararadhya, P.V. Ranjan, Efficient zone based feature extraction algorithm for handwritten numeral recognition of popular south Indian scripts, Journal of Theoretical and Applied Information Technology 7 (1) (2009) 1171–1180.

[23] N. Shanthi, K. Duraiswamy, A novel SVM-based handwritten Tamil character recognition, Pattern Analysis and Application, 2009, Published Online ⟨http://www.springerlink.com/content/k708558648652q12/⟩.

[24] T.V. Ashwin, P.S. Sastry, A font and size independent OCR system for printed Kannada documents using support vector machines, Saadhanaa 27 (Part 1) (2002) 35–58.

[25] B.B. Chaudhuri, U. Pal, Skew angle detection of digitized Indian script documents, IEEE Transactions on Pattern Analysis and Machine Intelligence 19 (2) (1997) 182–186.

[26] C. Luh Tan, A. Juntan, Digit recognition using neural networks, Malaysian Journal of Computer Science 17 (2) (2004) 40–54.

[27] M.B. Sukhswami, P. Seetharamulu, A. Pujari, Recognition of Telugu characters using neural networks, International Journal of Neural Systems 6 (3) (1995) 317–357.

[28] M. Wellner, J. Luan, C. Sylvestor, Recognition of Handwritten Digits using Neural Network, ⟨http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.136.9800⟩.

**About the Author**—APURVA A. DESAI received his Masters degree in Statistics, securing first rank in the South Gujarat University, Surat. He also received his Ph.D. from the same University in the year 1996. He has teaching and research experience of over 18 years in computer science. Presently he is working as a Professor of computer science in the university and also heading the department. He has many research papers and books to his credit. His major research areas include optical character recognition and data mining.