

# Recognition of Individual Handwritten Letters of the Farsi Language using a Decision Tree

Atefe Matin Niya  
Department of Computer Engineering  
Dezful Branch  
Islamic Azad University, Dezful, Iran

Hedieh Sajed  
Department Of Computer Engineering  
Amirkabir University of Technology,  
Tehran, Iran

## ABSTRACT

In this study, in order to recognize Farsi handwritten letters, first the pre-processing operation is done on the letters' images including normalization, thinning, reduction, noise reduction, etc., and then the feature vector of the letters is extracted using the first to the fourth momentums from the second level of wavelet transform and contourlet transform. A combination of decision-tree methods is used for the final recognition of letters.

The database used in this study is the "Hoda" handwritten letter collection. The mean recognition rate in this combinational method is 97.89%.

## Keywords

recognition of Farsi letters, handwritten recognition, pattern application, decision tree, wavelet transform, contourlet transform.

## 1. INTRODUCTION

Optical character recognition is an important, applied and very active branch in pattern recognition. The main problem in this branch of computer sciences is to recognize characters, subwords and words whose images are available. Solving this problem can have a tremendous effect on the effective relationship between man and machine and can also be a great help in the automation of written document processing. Generally, a document recognition system can be on-line or off-line. If time information on the writing sequence is not available, the system will be off-line and its data will be usually obtained by scanning the images of previously written texts. If the points' time sequence at the time of writing is available, on the other hand, the system will deal with on-line data.

This study aims to present a new method for the off-line recognition of Farsi individual letters.

To date, relatively fewer studies have been conducted on the off-line recognition of Farsi letters compared to Latin languages. A concise review of the most important studies in this field will be presented in section 2. Section 3 will examine the proposed method and its different subsections. Section 4 will present the experimental results of the proposed method and section 5 will express the conclusion and future solutions.

## 2. LITERATURE REVIEW

Recently, different approaches for writer identification have been proposed. A scientific validation of individuality of

handwriting is performed by Srihari et al. [1]. In this study handwriting samples of 1500 individuals, representative of the U.S. population with respect to gender, age, ethnic groups, etc., were obtained. The writer can be identified based on Macro features and Micro features that are extracted from handwritten documents. Said et al. [2] proposed a global approach based on multi-channel Gabor filtering, where each writer's handwriting is regarded as a different texture. Bensefia et al. [3] used local features based on graphemes extracted from segmentation of cursive handwriting. Then writer identification is performed by a textual based information retrieval model. Schomaker et al. [4] presented a new approach, using connected-component contours codebook and its probability-density function. Also combining connected-component contours with an independent edge-based orientation and curvature PDF yields very high correct identification rates. Schlapbach et al. [5] propose a HMM based approach for writer identification and verification. Bulacu [6] evaluated the performance of edge-based directional probability distributions as features in comparison to a number of non-angular features. Marti et al. [7] extracted a set of features from handwritten lines of text. The features extracted correspond to visible characteristics of the writing, for example, width, slant and height of the three main writing zones. In Ois et al. [8] a new feature vector is employed by means of morphologically processing the horizontal profiles of the words. Because of the lack of the standard database for writer identification, the comparison of the previous

## 3. THE PROPOSED METHOD FOR FARSI LETTER RECOGNITION

### 3.1 The pre-processing phase

In the pre-processing phase, images are first transformed into equal dimensions of  $24 \times 24$  pixels for normalization. Since the image lengths and heights may not change by the same ratio (i.e. they may stretch in one direction), images must be thinned. Figure (1) shows an example of thinning.



Fig 1: The thinning process

The noise caused by optical scanners leads to points such as stains, discrete line segments, connection between lines, filling of the existing holes in some letters' images, etc.

Also, different distortions should be considered including local changes, the curvature of letters' angles, letters' deformation or corrosion. Prior to the recognition stage, these defects should be fixed. The most important reason for noise reduction is to reduce error in the recognition phase. Also, noise reduction will reduce the image file size which in turn will reduce the time needed for future processing and storing. In this method, by determining a proper threshold level, we shall reduce noise and then transform grey images to black & white ones. This operation reduces the data size considerably.

In the next stage, images should be binarized so that an operation such as thinning is possible.

## 2.2 The representation phase (feature extraction)

After the pre-processing phase, properties must be extracted from images. The representation phase is a very important phase in OCR systems because the results obtained from this phase have direct impacts on the recognition phase quality. In the representation phase, each input pattern is assigned a feature vector or code which represents that pattern in the feature space and makes that pattern different from other patterns.

In this article, the feature extraction phase is divided into two parts:

The first part is feature extraction from wavelet transform coefficients: first, we will apply two levels of wavelet transform (DWT) on images using the Haar scaling function. The wavelet transform will cause images to cross high-pass and low-pass filter sets continuously so that data rate remains unchanged. Then, among the four obtained subbands (cA, cH, cD, cV), we will select the most effective one by means of the trial and error method. The selected subbands in this experiment are 3 subbands of cA, cH, cV that have the most similarities to the existing images. In the 2D wavelet transform, a 2D scaling function is required like  $\varphi(x,y)$  in which each one is the product of two 1D functions. Except for multiplications that produce 1D results like  $\varphi(x)\Psi(y)$ , the four remaining multiplications produce a separable scaling function:

$$\varphi(x,y) = \varphi(x) \varphi(y) \quad (2-1)$$

and separable "direction sensitive" wavelets:

$$\psi^H(x,y) = \psi(x)\varphi(y) \quad (2-2)$$

$$\psi^V(x,y) = \varphi(x)\psi(y) \quad (2-3)$$

$$\psi^D(x,y) = \psi(x)\psi(y) \quad (2-4)$$

These wavelets measure function changes (severity changes for images) in different directions:

$\psi^H$ : measures changes along columns (such as horizontal edges)

$\psi^V$ : respond to changes along rows (such as horizontal edges)

$\psi^D$ : corresponds to changes along diagonals.

### 2.2.1 The Haar transform

One of the operations related to imaging in multi-precision analyses is the Haar transform whose basic functions are the simplest and oldest of orthogonal wavelets. The Haar transform is expressed as matrix as follows:

$$T = HFHT \quad (2-5)$$

where F is a  $N \times N$  matrix, and H is the  $N \times N$  Haar transform. The matrix transpose is required because H is not symmetric. For Haar transform, H includes Haar basic functions, i.e.  $hk(z)$ . They were defined on the continuous range of  $z \in [0, 1]$  for  $k = 0, 1, 2, \dots, N-1$ , where  $N = 2n$ .

The second phase of feature extraction from the contourlet transform: A contourlet transform is a non-separable directional 2D transform that is used to describe subtle details and curves. Contourlet transforms efficiently describe smooth contourlets that are main and important components in natural images. Unlike other transforms which are created first in a continuous area and are then discretized for data sampling, a contourlet transform starts from a discretized area with the help of bank filters; then, it converges to a continuous area through an analytical multi-resolution framework. Like the first phase of images, we will obtain the contourlet transform coefficients in two levels and select the best subbands. To select the best subbands, we will use the trial and error method. In this phase, the best results were obtained from the union of four subbands.

Once the coefficients are obtained from the above transforms, image properties are extracted by the first to the fourth moments, i.e. mean, var, skewness, and kurtosis; subsequently, they are stored in a file.

Using the feature extraction method and a 2D wavelet transform, 36 properties are obtained, whereas using a contourlet transform, 40 properties are obtained; finally, a vector with 76 properties is formed. We obtain this vector for all of the images in the database and store them in a file. The group number is added to the end of the feature vector. This numbering is done based on the following table:

Table 1. numbering is done based

Lable	Persian Spell	Lable	Persian Spell
0	ا	18	ط
1	ب	19	ظ
2	بـ	20	ع
3	ت	21	غ
4	ث	22	ف
5	ج	23	ق
6	چ	24	ک
7	ح	25	گ
8	خ	26	ل
9	د	27	م
10	ذ	28	ن
11	ر	29	و
12	ز	30	ه
13	ژ	31	ی
14	س	32	نـ

15	ش	33	ا
16	ص	34	هـ OR هـ
17	ض	35	هـ

## 2.3 The recognition phase

In this phase, we will present a classification model for letter recognition based on the obtained feature vectors and Table 1.

### 2.3.1 The decision tree

Trees in artificial intelligence are used to represent different concepts such as sentence structures, equations, game modes, and so on.

Decision tree learning is a method for approximation of objective functions with discrete values. This method is resistant to data noise and can learn the seasonal composition of conjunctive statements. It is a popular inductive learning algorithm that has been used in several applications successfully.

A decision tree is a tree in which samples are classified in a way that they grow down from the root and finally reach leaf nodes.

In this article, the C4.5 decision tree is used for letter recognition.

### 2.3.2 The C4.5 decision tree

C4.5 builds decision trees from a set of training data in the same way as ID3 using the concept of information entropy.

The training data is a set  $S = s_1, s_2, \dots$  of already classified samples. Each sample  $s_i = x_1, x_2, \dots$  is a vector where  $x_1, x_2, \dots$  represent attributes or features of the sample. The training data is augmented with a vector  $C = c_1, c_2, \dots$  where  $c_1, c_2, \dots$  represent the class to which each sample belongs.

At each node of the tree, C4.5 chooses one attribute of the data that most effectively splits its set of samples into subsets enriched in one class or the other. Its criterion is the normalized information gain (difference in entropy) that results from choosing an attribute for splitting the data. The attribute with the highest normalized information gain is chosen to make the decision. The C4.5 algorithm then recurses on the smaller sublists.

This algorithm has a few base cases.

- All the samples in the list belong to the same class. When this happens, it simply creates a leaf node for the decision tree saying to choose that class.

- None of the features provide any information gain. In this case, C4.5 creates a decision node higher up the tree using the expected value of the class.

- Instance of previously-unseen class encountered. Again, C4.5 creates a decision node higher up the tree using the expected value.

## 4. IMPLEMENTATION AND EXPERIMENTAL RESULTS

MATLAB software has been used for implementation.

The "Hoda" data set is used as a training and experimental data set. This set includes 88351 samples, of which 70645 samples are used for training and 17706 are used for experimentation.

In the proposed method, the mean recognition rate of Farsi letters was obtained as 97.89%.

Table 2 shows the implementation results in the confusion matrix and table 3 shows the final results for each classes in data base.

## 5. SUMMARY, DISCUSSION AND CONCLUSION

This article presents a simple method for the recognition of Farsi letters based on recognition with the C4.5 recognition method. Thus, first the feature vector was extracted by the second level wavelet transform and the second degree contourlet transform; then, the image vectors were classified using these vectors. The pre-processing phase plays an important role in reaching desired results.

Based on the experimental results, it seems that the recognition of the letters' main bodies has proper results with the C4.5 method, but these tiny movements (including points, hats, etc.) eventually specify the letters' final shapes. Of course, this problem is more frequent in the groups where the differences between members lie in points, such as "پ، ب، " "ج، چ، ح، خ" and "ت، ث" groups. Despite problems in the recognition of tiny movements, it is probably better to use a more complicated decision tree in addition to the initial recognition.

Table2.confusion matrix

0	7	4	0	3	8	2	0	0	0	10	13	6	5	7	0	0	0	0	2	3	4	9	2027
4	11	1	0	2	1	3	0	0	0	0	2	0	1	5	1	0	0	1	6	16	18	1996	0
5	1	3	5	3	3	2	0	6	0	6	0	0	3	5	3	3	3	6	7	8	1994	7	9
15	8	5	1	1	1	1	0	1	1	1	0	0	4	5	0	0	0	1	25	1966	3	7	0
10	11	4	3	1	1	2	0	3	1	1	1	0	8	0	2	5	3	6	1721	12	0	4	1
6	0	9	4	3	4	0	0	3	0	2	0	0	0	0	27	34	36	1975	3	1	1	0	0
16	0	1	4	12	5	0	1	10	0	2	0	0	1	0	14	18	1894	19	2	0	1	0	0
8	2	12	2	1	2	2	0	0	0	3	0	0	2	0	33	1979	5	26	1	1	3	0	0
3	0	10	3	2	0	0	0	1	0	10	1	0	7	1	1934	3	2	9	4	0	3	0	5
3	2	1	2	0	0	0	0	1	0	9	6	2	15	2022	1	1	0	0	6	3	1	4	4
0	6	4	1	0	0	0	0	0	0	17	19	0	1765	10	2	0	0	0	13	7	7	0	2
0	0	0	0	0	0	0	0	0	0	0	9	2050	0	0	0	0	0	0	0	1	0	4	11
0	13	0	0	0	1	0	0	0	0	22	2004	3	1	1	0	0	0	0	1	0	2	5	4
0	3	9	1	0	0	2	0	4	0	1944	10	0	4	3	6	2	0	0	4	4	1	2	3
0	7	0	0	3	1	18	40	6	2013	0	0	0	0	0	0	0	0	0	6	0	0	0	0
8	7	2	1	6	8	20	6	2001	4	0	0	0	0	0	1	3	6	6	2	0	3	3	0
1	5	0	2	2	2	38	2000	7	25	0	0	0	0	0	0	0	3	1	3	2	1	0	0
5	5	2	4	6	6	1958	12	5	4	0	0	0	0	0	0	2	5	3	3	2	2	0	1
3	1	4	4	27	1991	1	0	7	0	1	0	0	0	0	1	1	0	0	0	1	4	0	0
4	4	3	3	1832	6	2	1	3	3	0	0	0	0	0	2	0	0	0	0	0	2	0	0
11	0	29	1961	3	2	2	0	0	0	3	0	0	1	1	4	12	4	4	4	5	3	0	3
2	0	1916	4	0	3	1	0	3	0	7	0	0	4	0	18	3	3	4	1	0	2	1	2
6	1955	2	2	2	1	1	1	3	1	1	1	2	1	0	2	0	1	1	4	5	3	6	0
1926	1	0	8	3	3	1	0	0	0	0	0	0	0	1	2	3	8	1	3	8	6	0	0
2	1	5	0	0	1	7	2	2	1	5	1	2	1	1	2	1	0	1	2	0	1	2	2

Table.3 The final results for each classes in data base

Results	character	Class #	Results	character	Class #	Results	character	Class #
97.66%	ک	24	98.76%	ز	12	95.93%	ا	0
96.59%	گ	25	97.60%	ژ	13	97.46%	ب	1
97.98%	ل	26	97.64%	س	14	97.35%	پ	2
96.56%	م	27	96.55%	ش	15	96.75%	ت	3
97.70%	ن	28	98.56%	ص	16	96.43%	ث	4
96.49%	و	29	98.04%	ض	17	95.39%	ج	5
96.29%	ه	30	98.32%	ط	18	94.89%	چ	6
96.95%	ی	31	99.12%	ظ	19	97.31%	ح	7
98.11%	نـ	32	98.51%	ع	20	95.58%	خ	8
99.23%	آ	33	97.51%	غ	21	98.42%	د	9
99.49%	هـ هـ	34	97.31%	ف	22	98.18%	ذ	10
99.71%	هـ	35	97.97%	ق	23	97.06%	ر	11

## 6. REFERENCES

- [1] F. Solimanpour, J. Sadri, C.Y. Suen, “**Standard databases for recognition of handwritten digits, numerical strings, legal amounts, letters and dates in Farsi language**”, in: Proceedings of the 10th International Workshop on Frontiers of Handwriting Recognition, La Baule, France, 2006, pp. 3–7
- [2] M. Hanmandlu, O.V. Ramana Murthy, Vamsi Krishna Madasu, “**Fuzzy Model based recognition of handwritten Hindi characters**”, Digital Image Computing Techniques and Applications, 0-7695-3067-2/07 © 2009 IEEE.
- [3] Abulhaiba S. H., Mahmood S. A., Green R. J., “**Recognition of handwritten cursive Arabic characters**”, IEEE Trans. On PAMI, Vol. 16, No. 6, pp. 664-671, 1994.
- [4] Altuwaijir M. and Bayoumi M., “**Arabic text recognition using neural networks**”, Proc. Int. Symp. On Circuit and Systems (ISCAS’94), pp. 415-418, 1994.
- [5] R. Plamondon; S. Srihari; “**On-line and Off-line Handwritten Recognition: A Comprehensive Survey**”; *IEEE Trans. On Pattern Analysis and Machine Intelligence*; Vol. 22; No. 1; January 2000; pp. 63-83.
- [6] N. Arica; F.T. Yarman-Vural; “**An overview of character recognition focused on off-line handwriting**”; *IEEE Trans. On Systems; Man and Cybernetics-Part C: Application and Reviews*; Vol. 31; No. 2; May 2009; pp. 216-233.
- [7] M. Ziaratban, K. Faez, F. Faradji, “**Language-based feature extraction using template-matching in Farsi/Arabic handwritten numeral recognition**”, in: Proceedings of the 9th International Conference on Document Analysis and Recognition, vol. 1, Curitiba, Brazil, 2007, pp. 297–301.
- [8] Teltscher, H. O. “**Handwriting-Revelation of Self**”, Hawthorn Books Inc. Publishers, New York, 1971.
- [9] J. C. Simon., “**Off-line cursive word recognition. Proc.**”. IEEE, 80(7):1150–1160, 1992.
- [10] Tal Steinherz, Nathan Intrator, and Ehud Rivlin. “**Skew detection via principal components analysis. In Proceedings of the 5th International Conference on Document Analysis and Recognition**”, pages 153–156. IEEE Computer Society
- [11] Ebrahimi.A. Kabir.E., “**A pictorial dictionary for printed Farsi sub words**”, Pattern Recognition Letters 29 (2008) 656–663.
- [12] R. Jacobs, “**Methods for combining expert’s probability assessments**”, Neural Computation, vol. 7, pp. 867-888, 1995.