

Feature Extraction Techniques for Handwritten Text in Various Scripts: a Survey

Vandita Singh, Bhupendra Kumar, Tushar Patnaik

Abstract—Optical Character Recognition (OCR) Systems aim to recognize text and bring it to editable form from the given document image, where the input text can be in machine printed, hand written or hand printed form. Many recognition systems have been developed for languages based on various scripts and digits all over the world, taking input in either of the online and offline modes, with considerable efficiencies. These systems have proved to be highly applicable in the fields of Banking, Education, IT systems and Postal Sector for digitization of processes and automated information retrieval. In this paper, we present a survey of techniques for recognition of handwritten and hand printed documents in off-line mode, with an emphasis on the Feature Extraction phase and the corresponding classification technique has also been mentioned with the recognition rates achieved.

Keywords— Optical Character Recognition, Feature Extraction, Classification.

I. INTRODUCTION

The domain of Digital Image Processing and computer vision give a way for devising an approach for development of an optical character recognition system meant to understand the digital image and infer from it the editable document required for various applications such as automated data entry and retrieval, bank cheque analysis, hand written pictogram interpretation, hand written formulae retrieval, address retrieval and verification and many more. Various researchers all over the world have proposed and applied methodologies to serve solutions to these problems. For many languages and modes of input for all these applications, the OCR system holds to be robust. However more efficient systems may also be developed, keeping in view the requirements and nature of input documents for certain specific applications such as for OMR forms processing, where hand printed data can be recognized by machines in addition to the recognition of encircled bubbles.

A. Background

Background of the research work carried out in field of Hand written document image recognition goes back to decades. The history of OCR in field of Handwritten Text Recognition has been discussed by J.R Prasad and Kulkarni [1] as trends in Handwriting recognition. The earliest work in this domain according to the mentioned survey work [1], has been reported to be done in the early sixties and seventies. Character images were then modelled as being composed of blocks defined geometrically by the co-ordinates of their vertices and the specification of edge information. Later integrated segmentation and interpretation systems came into being.

Then there came into existence, the algorithmic and computational techniques to process the image at high level (edge finding, region growing and segmentation). Further improvements were made to these techniques to obtain more efficient algorithms for classification such as Hidden Markov Models (HMM), Dynamic Programming, Neural networks etc.[1]. Thus the trend moved on towards devising automated techniques for recognition, incorporating efficient algorithms for pre processing, feature extraction and classification. Also use of combination of recognizers, use of lexicons, dictionaries and language models [1] is being introduced to achieve higher recognition rates.

B. Handwritten v/s Hand printed Text

Certain terms such as handwritten, hand printed, on-line and off-line, often misinterpreted, need to be clarified before moving further. The term Hand printed document image implies that input image comprises of mono-spaced characters, i.e. each character as it's filled out is the same distance apart as all the other characters whereas in handwriting we have characters that connect, the extreme form being cursive whereas Hand printed input document has more or less uniform height or width. Thus, as many variable elements are not introduced as in straight handwriting (unconstrained). This may lead to a minimisation of variability. Hand printed document images also imply that there must be present uniform base-line character images (same horizontal base-line). The other terms being off-line and on-line recognition [2] modes of obtaining input for recognition can be explained as follows. On-line document image recognition implies to storing in order the two dimensional co-ordinates of the successive points of the writing as a function of time, thus the spatio-temporal information such as order of strokes made by the writer, information about pressure and angle of the pen is readily available. In case of off-line document where acquisition is done prior to recognition, the image consists of the complete data to be written. Here the spatio-luminance of the image is analysed. Therefore, more challenges would be present while recognizing documents in offline mode since we have only static information about the document.

II. FEATURE EXTRACTION-A CRUCIAL PHASE IN THE OCR PROCESS

When the pre-processing and the desired level of segmentation (line, word, character or symbol) has been achieved, we apply some feature extraction technique to the segments to obtain features, which is followed by application of classification and post processing techniques. It is essential to focus on the feature extraction phase as it has an observable impact on the efficiency of the recognition system. Anil K. Jain, Taxt [3] emphasise on the fact that feature selection of a feature extraction method is the single

Manuscript received Feb 02, 2013.

Vandita Singh, C-DAC-Noida, Noida, India.

Bhupendra Kumar, C-DAC-Noida, Noida, India.

Tushar Patnaik, C-DAC-Noida, Noida, India.

most important factor in achieving high recognition performance [3]. According to Devijver and Kittler, (page 19 of [3]), definition of feature extraction has been given as “extracting from the raw data information that is most suitable for classification purposes, while minimising the within class pattern variability and enhancing the between class pattern variability”. Thus, selection of a suitable feature extraction technique according to the input to be applied needs to be done with utmost care. Also, efficiency of classification and therefore, that of recognition depends upon the feature extraction phase. As we obtain an n-dimensional feature vector from this phase that is fed into the classifier for further processing, feature set hence determined needs to be optimal enough to give efficient results, i.e. it must take into consideration the curse of dimensionality [3] as mentioned in the survey work by A.K Jain and Taxt. The extracted features must also be invariant to the expected distortions and variations that a character may undergo [3].

Taking into consideration all these factors, it becomes essential to look at the various available techniques for feature extraction in a given domain, covering vast possibilities of cases. Therefore we move forward to have a look over the techniques developed till date for handwritten and hand printed character image recognition in off-line mode, where we cover various languages and scripts- Assamese, Persian, Thai, Devanagari, Chinese and Roman characters and digits.

III. FEATURE EXTRACTION TECHNIQUES

The basic classification can be done mainly on the basis of representation of input document as on-line and off-line, where in each category, we may have numerous techniques- statistical, structural, rule based etc. for both hand printed and handwritten document recognition in various language scripts. The various techniques can be enumerated as follows-

A. Hand printed Numeral Recognition

Various scripts discussed for numeral recognition are- Roman, Assamese and Persian. One of the techniques proposed by J.R Parker [4] to recognize hand printed digits is to use scalable vector templates, which generate templates with same scale and line width attributes as an arbitrary image[4]. There was the problem to match ‘1’ or any other object that had extreme height to width ratio. Shridhar et al [5], had used a combination of global-left and right profiles of external contours, difference of left and right profiles and local features- width of character, ratio, location of maximum and minimum and right, left peaks of the digit. A syntactic recognition algorithm was implemented by formulation of production rules [5], resulting into overall accuracy as 99%. Heutte et al [6], used seven features-both statistical and structural, combined to represent the numerals. These were intersection with straight lines, invariant moments, holes and concave arcs, extrema, end points and junctions, profiles, projections [6]. Statistical classifier [6] was employed and 90.8% accuracy was achieved. Liu et al [7] made use of variants of directional features- chain code, gradient feature with Sobel and Kirsh operators, and peripheral direction contributivity (PDC)[8] in complementation with profiles structural features[9] and concavities structural features[10] to compose ten feature sets[7]. Classifier system implemented by the authors[7], composed of eight classifiers (k-NN, MLP, RBF, PC, LVQ, DLQDF, SVC-poly, and SVC-RBF)[7] to give the highest accuracy as 99.58% when

using SVC with RBF kernel (SVC-RBF)[7]. Koerich [11] have also made use of profiles as complementary features for recognition of hand-printed pattern.

G.Siva Reddy et al.[12] achieved remarkable recognition rates 97.6% for offline Assamese handwritten numerals using VPP-HPP, Zonal DCT, CCH, Pixel level features and a combination of all the four features[12] to represent features and vector quantization[12] for modelling[12]. Omid Rashnoodi et al.[13] have performed offline recognition of handwritten Persian digits using five feature sets, composed of-variance, co-variance, central moments, median and number of pixels per each square[13] and two classifiers-SVM[13] and k-Nearest Neighbour[13] to give 91.3% and 92% as recognition rates, respectively.

B. Chinese and Thai Characters

Another research methodology for feature extraction in Chinese hand printed characters was to extract features using statistical methods. Dominant Point Method [14] was the most popular for the said problem, for which several algorithms had already been developed. However, the authors [14] felt a need to detect only the points with a very sharp curvature. The authors had followed the Rosenfeld-Johnson algorithm[14], whose basic concept had been stated as “to calculate curvature of each point in line, then the points with local maximum in curvature are designated as dominant points”[14]. The authors had formulated six equations to denote the six primitives. The recognition rate [14] achieved was 84.45%. This methodology worked on the need to break traced lines into segments for easier recognition. This algorithm was affected by the presence of undesirable dominant points due to irregularities of lines; however authors had proposed and applied re-merging of line segments which resulted into success. Another technique for feature extraction was applied on Hand printed Thai characters [15], which combined both global and local features of the characters. Statistical Features were represented by global features to define the shape of the characters, which was further used to group similar shape characters together. Local features were used to represent the symbolic structures for the characters. The structural features identified were loops, end points, junction points and curls. Authors had improved efficiency for recognition of unconstrained (handwritten characters). Also this methodology resulted out to be robust and an improvement over recognition of structural feature extraction, with efficiency of 87.32%[15] for the combination of global and local features.

C. Devanagari Script

There has been considerable development in field of Indian language OCR, where recognition of Hand-printed characters in Devanagari script has been done by implementing a combination of features to yield efficient results. Features implemented and compared [16] in this paper are as follows-

- Kirsch Directional Edges
- Distance Transform
- Chain Codes (Freeman Chain Codes)
- Directional Distance Distribution
- Neighbourhood Pixels Weight
- Total Distances in four Directions

The classifiers used were MLP and SVM [16], that was fed with multiple features and results were compared. It was noted that lowest efficiency was obtained by using Kirsch

Directional edges and highest using Gradient. Overall recognition rate obtained was 94.3% using SVM (Support Vector Machines) [16].

Satish Kumar [17] has also thrown some light on the analysis of following mentioned features, in addition to the five features mentioned according to [16]. These features [17] are-

- Zoning
- Profiles
- Chamfer
- Histogram
- Crossings
- Chessboard
- Chain code Histogram
- Gradient (Sobel)
- Total Distances in 4 directions combined with Gradient-TdistGrd-200
- Neighbouring pixels Weight combined with Gradient-NpwGrd-200.

Evaluation has been done using three classifiers-SVM(RBF),k-NN and MLP[17].Best result 93.5% has been stated as obtained using Gradient (Sobel) -200 [17]with SVM as classifier.

Various feature extraction techniques for hand printed Devanagari numeral and character recognition have been enlisted by R. Jayadevan et al [18], where efficiencies of each of the techniques have been compared. Likewise a wide variety of features that have been used for Devanagari numerals and characters both, have been covered by the authors [18], some of them are as follows-

- Structural
- Statistical
- Density features
- Moment features of right, left, upper, and Lower profile curves
- Descriptive component features
- Shape descriptors.

A comparative study has been done by the authors and feature extraction techniques with corresponding recognition results [18] of each of the techniques have been tabulated. Ashutosh Aggarwal, Rajneesh et al [19] had worked on recognition of offline Devanagari characters, obtaining gradient feature vectors and applying SVM (RBF) as the classifier to yield recognition rates as 94%.

D. English language Hand printed and Handwritten Characters

Gilewski et al [20] performed in-depth laboratory work to study techniques for Handwriting Recognition. The authors had applied the fuzzy feature extraction algorithm [20], with the main aim being to extract the defined segments from the character image. Multilayer feed forward Neural Network was used in the laboratory in order to realize the proposed algorithm. Cursive Character Recognition using Multiple Feature Extraction Techniques was applied by Rafael M.O. Cruz et al[21] while making use of an ensemble classifier for recognition. Nine feature sets were constructed. Other techniques [21] can be enlisted as follows-

- Structural Characteristics, where horizontal, vertical and radial histograms were computed.
- Modified Edge Maps
- Image projections

- Multi Zoning
- Concavities Measurement
- MAT-Based Gradient directional features, where Medial Axial Transform was applied to binary input to transform it into pseudo-greyscale image.
- Gradient Directional Features
- Median Gradient Features
- Camastra 34D Features

These nine feature sets [21] were given to the ensemble system of classifiers. The motivation authors [21] felt for building ensemble system were as follows- Errors made by classifiers with different feature extraction methods do not overlap and Divide and Conquer paradigm, i.e. using each feature set separately and combining their results would prove to be more efficient. Authors have presented a comparison of results obtained for each feature set and each classifier. However, the new technique of Modified Edge Map produced the best overall results, 84.37% [21].

Offline Cursive Character recognition was performed by Anshul Gupta et al [22] incorporating heuristics based segmentation algorithm. The recognition was performed at the word level which required segmentation of word into independent characters (uppercase and lowercase).Fourier Descriptors [22] were used in order to extract features, namely Fourier Magnitude [22] and Fourier Angle [22].Three networks were considered as classifiers-Multi-layer perceptron (MLP), Radial Basis Function (RBF) and Support Vector machine (SVM).Ten variations of 26 word images were considered. The proposed system [22] achieved recognition accuracy of 98.74% on the training dataset using SVM as the classifier .As stated by the authors; it had outperformed RBF and MLP [22]. The summarization of the above mentioned feature extraction and classification techniques with recognition rates achieved in each case have been tabulated in Table I.

Table I. Feature Extraction techniques and corresponding classifiers

S. No.	Languages	Feature Extraction Technique	Classification Technique	Result
I.	Hand printed Numerals(Roman digits) [4]	Vector Template Representation	Template Matching	94.3%(Average result)
	Persian Digits [13]	Five Features' set	SVM	91.3%
			k-NN	92%
II.	Chinese Characters[14]	Freeman chain codes	Statistical Discriminant Analysis	84.45%
III.	Thai Characters[15]	Local-symbolic and Global-statistical	NLP Neural network Models	87.32%
IV.	Devanagari Characters [16], [17]	TDIST and four other techniques	MLP and SVM	94.3% (Overall using SVM)
		Gradient(Sobel)	SVM(RBF)	93.5%
V.	English language(Roman Script) [21], [22]	Nine features' set	Ensemble System	84.37% (Modified Edge Maps as Features)
		Fourier Descriptors	SVM	98.37%

IV. CONCLUSION

The intent of covering these techniques was to have an overview of the existing feature extraction techniques and perform a comparative analysis of various such techniques in relation with the nature of input. Much of work has been done taking into account statistical and structural features, both independently and in combination as well. However, more work is required to be done using rule based methods for feature extraction, which may lead to even better results in future. It can be seen that use of profiles has been prevalent especially for recognition of roman digits. Statistical techniques such as multi zoning and edge maps yielded very good efficiencies for hand written Roman alphabet recognition, for which the reason could be a large feature set is obtainable in these cases. For classification, the best efficiencies have been obtained by using intelligent systems, specifically SVM as the classifier. However the ensemble system also holds out to be a strong competitor for being chosen as the classifier. The techniques for feature extraction and classification hence discussed would prove to be helpful to develop a clear understanding of the concepts essential to build a recognition system for offline hand written documents, even applicable to scripts other than those discussed in this survey.

REFERENCES

- [1] Jayashree R. Prasad, U.V. Kulkarni 2010. "Trends in Handwriting Recognition" IEEE 2010
- [2] Plamondon and Srihari, "On-line and Off-line Handwriting Recognition-A Comprehensive Survey" IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol 22, No.1, January 2000
- [3] Ivind Due Trier, Anil K. Jain and Torfinn Taxt, "Feature Extraction Methods for Character Recognition-a Survey", Pattern Recognition Vol.29, No.4, p.p.641-662 (1996)
- [4] J. R. Parker, "Vector Templates and Hand printed Digit Recognition", 1051-4651, IEEE-1994
- [5] M. Shridhar and A. Badreldin, "A High Accuracy Syntactic Recognition Algorithm for handwritten Numerals", Proceedings of the IEEE International Conference on Systems, Man and Cybernetics, Vol. SMC-15, No. 1, January/February 1985
- [6] L. Heutte, T. Paquet, J. Moreau, Y. Lecourtier and C. Olivier, "Combining Structural and Statistical Features for the Recognition of Handwritten Characters", ICPR, Vol. 19, pp. 629-641, 1998.
- [7] C.-L. Liu, K. Nakashima, H. Sako and H. Fujisawa, "Handwritten Digit Recognition: Benchmarking of State-of-the-Art", Pattern Recognition, No. 36, pp. 2271- 2285, 2003.
- [8] A. de S. Britto Jr., et al., "Improvement in handwritten numeral string recognition by slant normalization and contextual information". Proceedings of the Seventh International Workshop on Frontiers of Handwriting Recognition, Amsterdam, 2000, pp. 323-332.
- [9] J.T. Favata, G. Srikantan, S.N. Srihari, "Hand printed character/digit recognition using a multiple feature/resolution philosophy". Proceedings of the Fourth International Workshop on Frontiers of Handwriting Recognition, Taipei, 1994, pp.57-66.
- [10] J.T. Favata, G. Srikantan, S.N. Srihari, "Hand printed character/digit recognition using a multiple feature/resolution philosophy", Proceedings of the Fourth International Workshop on Frontiers of Handwriting Recognition, Taipei, 1994, pp.57-66.
- [11] L. Koerich, "Large Vocabulary Off-line Handwritten Word Recognition", Ph. D. Thesis, Ecole de Technologie Superieure, Montreal - Canada, 2004.
- [12] G.Siva Reddy, Puspanjali Sharma, S.R.M Prassana, C.Mahanta, L.N. Sharma, "Combined online and Offline Assamese Numeral Recognizer" /978-1-4673-0816-8/12 IEEE 2012
- [13] Omid Rashnoodi, Asgher Rashnoodi and Aref Rashnoodi, "Offline Recognition of Handwritten Persian Digits using Statistical Concepts". International Journal of Computer Applications (0975 – 8887) Volume 53– No.8, September 2012
- [14] A. Amin and S. Singh, "Machine Recognition of Hand Printed Chinese Characters" Intelligent Data Analysis (1997) 101- 118

- [15] Ithipian Methasate, Sanparith Marukatat, Sutat Saetang and Thanarak Theeramunkong, "The Feature Combination Technique for Offline Thai Character Recognition System" IEEE, Eighth ICDAR 2005
- [16] S.Kumar, "Performance Comparison of Features on Devanagari Hand Printed Dataset", International Journal of Recent Trends in Engineering, Vol.1, No.2, May 2009
- [17] S.Kumar, "Study of Features for Handprinted Recognition", World Academy of Science, Engineering and Technology 60 2011
- [18] R.Jayadevan, Satish R.Kolhe, Pradeep M. Patil, Umapada Pal. 2011. "Offline Recognition of Devanagari Script-A Survey", IEEE transactions on Systems, Man and Cybernetics-Part C: Applications and Reviews, Vol.41, No.6, November 2011
- [19] Ashutosh Aggarwal, Rajneesh Rani, RenuDhir, "Handwritten Devanagari Character recognition Using Gradient Features", IJARCSSE Volume 2, Issue 5, May 2012
- [20] J. Gilewski, Phil Phillips, S. Yanushkevich, D.Popel, "Education Aspects-Handwriting Recognition using Neural Networks, fuzzy Logic" Pattern recognition and Information processing, vol 1,(1997) p.p. 39-47.
- [21] Rafael M.O Cruz, George D.C Cavalcanti, Tsang Ing Ren, "Ensemble Classifier for Offline Cursive Character Recognition Using Multiple Feature Extraction Techniques". IEEE
- [22] A.Gupta, M. Srivastava, C.Mahanta.2011, "Offline Handwritten Character Recognition" ICCAIE (2011)



Ms. Vandita Singh is currently pursuing M.Tech (Fourth semester) in Information Technology from C-DAC Noida. Her interest areas are Digital Image Processing, Database Management Systems and Software Engineering.



Mr. Bhupendra Kumar (Senior Technical Officer) joined CDAC in 2005, he received his M.Tech from IIT Allahabad with the specialization in wireless communication and computing. His interest areas are Advanced Image processing, pattern recognition, computer network, wireless network, MANETs. Currently he is involved in project "Development of Robust Document Image Understanding System for Documents in Indian Scripts".



Mr. Tushar Patnaik (Sr. Lecturer/Sr. Project Engineer) joined CDAC in 1998. He has eleven years of teaching experience. His interest areas are Computer Graphics, Multimedia and Database Management System and Pattern Recognition. At present he is leading the consortium based project "Development of Robust Document Image Understanding System for Documents in Indian Scripts".