

A Survey on Gujarati Handwritten OCR using Morphological Analysis

Vaidehi Patel, Prof. Abhinay Pandya

Department of Computer Engineering, LDRP Institute of Technology and Research, Gandhinagar, Gujarat, India

ABSTRACT

Handwriting recognition has been one of the active and challenging research areas in the field of image processing and pattern recognition. Handwriting recognition can be differentiated into two categories i.e. online handwriting recognition and Offline handwriting recognition. I mainly focus on offline handwriting recognition. It is also observed that majority of work for this language is for printed form rather than handwritten form. One can obtain very less OCR related research work for Gujarati script, especially for handwritten form. I use Freeman chain code for feature extraction and then used Hidden Markov method. There is no benchmark dataset is available so it is generated by collecting sample from more than 100 writers of different age group and gender.

Keywords: Handwritten Recognition, Freeman Chain Code, Pre-processing, Classification

I. INTRODUCTION

The form of the text is varied from the scanned document to the typed text in various fonts. For humans, recognition of text is a trivial task, but to make a machine that recognizes characters is extremely difficult. The current study is being focused on exploration of possible methods to develop an OCR system for Gujarati language when noise is present in the signal. In order to understand the core challenges, a thorough analysis of Gujarati Writing System has been done. To know and understand the research in this field, existing OCR systems are studied. The Prominence was on finding workable segmentation technique and diacritic handling for Gujarati words, and building a recognition module for these ligatures.

To develop an OCR system Gujarati text, a complete procedure is proposed as well as a testing application is also made. Test results from this are reported and compared with the previous work done in this area.

II. METHODS AND MATERIAL

A. Features of Gujarati Language

Gujarati is regional language of state Gujarat in India.

Gujarati is name of script which is written and spoken by people in Gujarat. Gujarati script has 34+2 constants and 11 vowels. The structures of some Gujarati characters are very similar to Devnagri script. Devnagri script has shirolekha on top of the characters, but Gujarati script has not. Gujarati script does not have the distinction of Lower and Upper Cases like English script. Gujarati script has combination of constants and vowels. Every vowel has a unique symbol, called vowels modifiers. Table 1 shows Gujarati vowels, consonants, matras, other symbols and some of the conjuncts [5].

Consonants
ક ખ ગ ઘ ડ ય છ જ ઝ ઞ ઠ ડ ઢ ણ ત થ દ ધ ન પ ફ ભ બ મ ય ર લ વ શ ષ સ હ ળ ક્ષ જ્ઞ
Vowels
અ આ ઇ ઈ ઉ ઊ એ ઐ ઓ ઘો ઘૌ ઘં
Some Conjuncts
ક જ્ય કલ ચ્છ દ્ઢ ત્ર સ્પ સ્ત વ્ વ્ઙ ત્ય ત્ર્ય ત્ત દ્
Consonant – vowel
જા ગી જી બુ બૂ હે કૃ કે પૌ ડો કૌ કં
Conjunct -vowel
જ્ઞા કે ક્લે ચ્છે ક્ષી ત્રુ સ્નૂ
Vowels modifiers
। િ ી ુ ૂ ૃ ૄ ૅ ૆ ે ૈ ૉ ૊ ો ૌ ્

Figure 1 : Gujarati Script[5]

B. Dataset Generation

There is no benchmark dataset is available so it is generated by collecting sample from more than 100 people of different age group and gender.



Figure 2 : Handwritten Characters

C. Feature Extraction

The objective of feature extraction is to capture the essential characteristics of the symbols, and it is generally accepted that this is one of the most difficult problems of pattern recognition. There are below techniques used.

1. Zoning
2. Projection Histogram Feature
3. Moments
4. Fourier Transform
5. Freeman chain code

Table 1 - Comparison of different feature extraction techniques

Techniques	Advantages	Disadvantages
Zoning	Simple, Easy to implement	Density of object pixel in each zone is calculated. Efficiency is low
Moments invariants	Easy recognize pattern field	Higher order moments are sensitive to noise and variation of writing style
Freeman Chain code	Process time small, Storage is lossless, easy to represent easy to implement	Length of FCC depend on starting point, may be revisit same node
Fourier Transform	Give valuable info about character structure. Recognize position shifted character	Not give accurate result. Tough to implement
Projection Histogram	Original histogram can be recovered, easy to implement	It is indiscriminate. It may increase the contrast of background noise while decreasing the usable signal

III. RESULTS AND DISCUSSION

CLASSIFICATION

The Extracted features are given as the input to the Classification process. A bag-of-key point extracted from the feature extraction approaches are used for classification.

1. Template matching
2. Neural Networks
3. Support Vector Machine (SVM)
4. K-nn Classification
5. Hidden Markov Model

See all of these techniques in details.

Template matching is simple technique of character recognition; depend on matching the stored templates with the character or word to be recognized. The matching operation finds out the similarity between two vectors. An input image is matched with set of already stored templates. The recognition rate of template matching is proportional to noise and image deformation. [9]

Neural Networks is composed of interconnected nodes that are connected via links. Learning is provided by example via training, or exposure to a set of input/output data (patterns), where the training algorithm adjusts the link weights.

A simple neural network
input layer hidden layer output layer

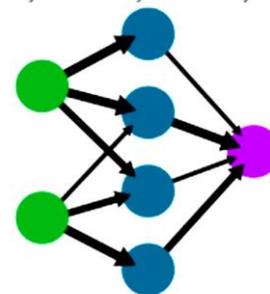


Figure 3 : Neural network

K-nn Classifiers is a nonparametric method used for classification. It is a Statistical method. So, basically the k-nearest neighbor algorithm (k-NN) is a method for classifying objects based on closest training examples in the feature space [4].

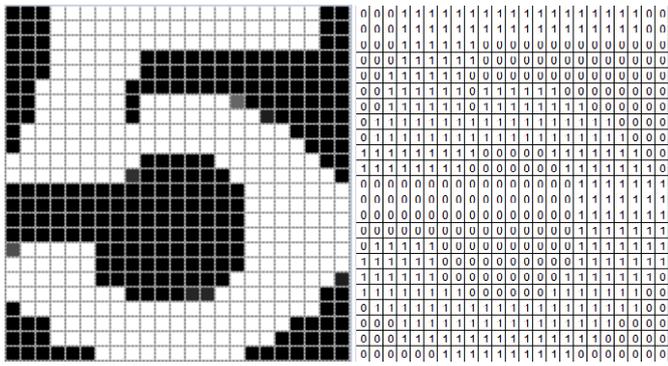


Figure4: The kNN classifier that is used here uses Euclidean distance method

Support vector machines (SVM), when applied to text classification provide high accuracy, but poor recall. One means of customizing SVMs to improve recall, is to adjust the threshold associated with an SVM. SVMs have achieved excellent recognition results in various pattern recognition applications.

Hidden Markov Model is a finite set of states, each of which is associated with a probability distribution. Transitions among the states are governed by a set of probabilities called transition probabilities. In a particular state an outcome or observation can be generated, according to the associated probability distribution. The probabilities for each candidate character are calculated. Then, the probabilities are counted to obtain a final best character-list for character recognition [1].

Table2: Comparison Table

Method	Advantages	Disadvantages
Template Matching	High Speed. Simple to implement No need to extract Features	Not effective when there are font discrepancy, font slant, font defilement. The method is not invariant to changes in illumination
Support Vector Machine	Good for small category dataset Easy to compute. Less time	High Complexity of training and execution. If dataset is big then results large number of SVMs which required more storage and time
Neural Networks	Higher recognizing ratio	Take more time for training

	with more training able to adapt to changes in the input data Accuracy is Greater	
K-NN algorithm	Simple Non parametric	Each step of distance is calculated so takes more time
Hidden Markov Model	Freedom to manipulate training and verification. The probability of observing sequence model and computed for each word.	Not completely automatic

IV. CONCLUSION

In this paper many techniques explained for the scanned document. Different OCR methods relevant to the survey of many papers. For developing faster OCR than Hidden Markov model method is very suitable because it is probabilistic model and not requires more training.

V. REFERENCES

- [1] Jinhong K. Guo and Matthew Y. Ma "Separating Handwritten Material from Machine Printed Text Using Hidden Markov Models" Published by IEEE 2001.
- [2] Sandhya Arora, Debotosh Bhattacharjee, Mita Nasipuri, Dipak Kumar Basu*, Mahantapas Kundu, 'Combining Multiple Feature Extraction Techniques for Handwritten Devnagari Character Recognition' published by IEEE 2008.
- [3] Debi Prasad Bhattacharya, Susmita Koner "English alphabet recognition using chain code and LCS" published by Indian Journal of Computer Science and Engineering (IJCSSE) ISSN : 0976-5166 Vol. 3 No. 2 Apr-May 2012.
- [4] Chhaya Patel and Apurva Desai , "Gujarati Handwritten Character Recognition Using Hybrid Method Based On Binary Tree-Classifer And K Nearest Neighbour" published by International Journal of Engineering Research & Technology (IJERT), ISSN: 2278-0181 Vol. 2 Issue 6, June - 2013

- [5] Dholakia, J.; Yajnik, A.; Negi, A., "Wavelet Feature Based Confusion Character Sets for Gujarati Script," Conference on Computational Intelligence and Multimedia Applications. International Conference on, vol.2, no., pp.366,370, 13-15 Dec. 2007
- [6] Sandhya Arora, Debotosh Bhattacharjee, Mita Nasipuri, Dipak Kumar Basu*, Mahantapas Kundu,'Combining Multiple Feature Extraction Techniques for Handwritten Devnagari Character Recognition' 978-1-4244-2806-9/08/\$25.00© 2008 IEEE
- [7] Anitha Mary M.O. Chacko and P.M. Dhanya 'A Comparative Study of Different Feature Extraction Techniques for Offline Malayalam Character Recognition '© Springer India 2015
- [8] Chaudhari Shailesh A., and Ravi M. Gulati. "An OCR for separation and identification of mixed English—Gujarati digits using kNN classifier."Intelligent Systems and Signal Processing (ISSP), 2013 International Conference on. IEEE, 2013.
- [9] Prasad, J.R.; Kulkarni, U.V.; Prasad, R.S., "Template Matching Algorithm for Gujrati Character Recognition," Emerging Trends in Engineering and Technology (ICETET), 2009 2nd International Conference on , vol., no., pp.263,268, 16-18 Dec. 2009