

REVIEW ARTICLE

OFFLINE HANDWRITING RECOGNITION USING GENETIC ALGORITHM

*Snehal S. Patwardhan and Dr. Deshmukh, R.R.

Department of CS and IT, Dr. B.A.M. University, Aurangabad, Maharashtra, India

Accepted 14th August, 2015; Published Online 30th September, 2015

ABSTRACT

The optical character recognition (OCR) is known to be one of the earliest applications of artificial intelligence. Genetic algorithm, which partially emulate human thinking in the domain of artificial intelligence, has been used in this study for OCR. Dealing with hundreds of features in character recognition systems is not unusual. This large number of features leads to the increase of computational workload of recognition process. There have been many methods which try to remove unnecessary or redundant features and reduce feature dimensionality. Experimental results show that application of genetic algorithms (GA) to features selection in a Marathi OCR results in lower computational complexity and enhanced recognition rate.

Key Words: Genetic Algorithms, Optical Character Recognition.

INTRODUCTION

Any pattern recognition system typically consists of a section which defines and extracts useful features from a pattern and uses a classifier to classify input patterns into different classes. OCR is the acronym for Optical Character Recognition. It has become an important and widely used technology for pattern recognition. OCR mechanism converts images with text into text documents using automated computer algorithms. In traditional recognition technique, images can be processed individually (.jpg, .png, and .gif files) or in multi-page PDF documents (.pdf) however OCR has even advanced into a newer field - Handwritten Recognition like genetic algorithms, which is also based on the simplicity of Character Recognition. (Pal and Chaudhari, 2004) Artificial intelligence elements like, artificial neural networks, genetic algorithms, fuzzy logic, expert systems, SVM etc. are tending to emulate the human brain and are one of the main interests of the researches nowadays.

The genetic algorithms were first suggested in 1975, by John Holland. Recently and currently using in a range of problems together with scheduling, images creating, planning strategy, predicting with dynamical systems, classification etc. Subhra provided the construction of genetic algorithm based neural network for parameter estimation of Fast Breeder Test Reactor (FBTR) Subsystem. Populations, encoding, crossover and mutation operators, static fitness functions, selection and termination are generally used to evaluate the performance of the Genetic algorithms. Genetic algorithms offer a particularly attractive approach for this kind of problems since they are generally quite effective for rapid global search of large, nonlinear and poorly understood spaces. Moreover, genetic algorithms are very effective in solving large-scale problems. (Pal *et al.*, 2007) Rest of the paper is organized as follows. In Section II Optical Character Recognition is discussed.

In Section III Application of Character Recognition is described. In section IV the Limitation of Character Recognition is described. In section V the Genetic algorithms are described with its working principles. In section VI Genetic Algorithm Based Classification And Recognition Techniques are described. In Section VII. Finally section VII is based on the Conclusion.

Optical character recognition

Fig. 1 shows the block diagram of an OCR system. The system involves 5 stages: Preprocessing, Segmentation, Feature extraction, Classification and Post processing. A typical OCR system may not include some of these stages. The recognition process starts by acquiring a digitized image. In the first stage the preprocessing of the image takes place. Recognition accuracy of the system depends on the image quality and amount of noise that exists in the image. All of the processes that improve image quality and prepare it for next stages are called preprocessing. Binarization, noise detection, smoothing and thinning are examples of these processes. Segmentation is the most important part of OCR systems.

There are two kinds of segmentation: first, distinguishing different components of a script like paragraphs, sentences and words; and second segmentation of a word to its characters. After segmentation a set of features is required for each character. In feature extraction stage every character is assigned a feature vector to identify it. This vector is used to distinguish the character from other characters. Some of most common approaches for feature extraction are Hugh transforms, moments and characteristic loci. The step, following the segmentation and extraction of appropriate features, is the classification and recognition of characters. Many types of classifiers are applicable to the OCR problem, among which, neural classifiers and distance function achieve very good results. Finally post processing stage tries to improve recognition results using additional information like word dictionary. (Soryani and Rafat, 2008)

*Corresponding author: Snehal S. Patwardhan

Department of CS and IT, Dr. B.A.M. University, Aurangabad, Maharashtra, India.

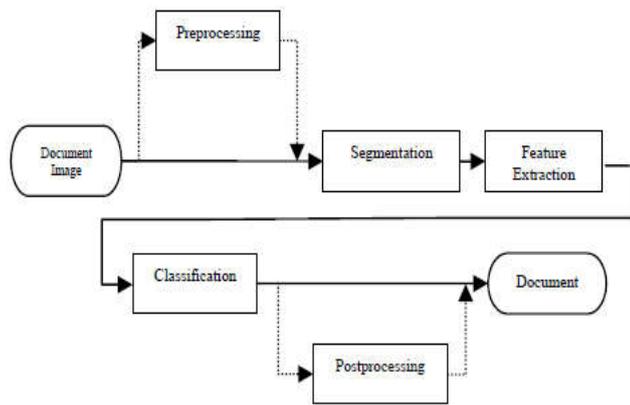


Fig. 1. Block diagram of an OCR system

Application of Character Recognition

OCR has enabled scanned documents to become more than just image files, turning into fully searchable documents with text content that is recognized by computers. With the help of OCR, people no longer need to manually retype important documents when entering them into electronic databases. Instead, OCR extracts relevant information and enters it automatically. The result is accurate, efficient information processing in less time. The character recognition systems can be classified into 2 categories.

Task Specific Readers

The task specific readers handle only specific document types. Some of the most common task specific readers read bank check, letter mail, or credit card slips. This considerably reduces the image processing and text recognition time. Some application areas to which task specific readers have been applied are:

Form Readers

There are various form readers available.

- Legal Form Readers
- Healthcare Patient Forms Readers
- Signature Verification
- Stamp Recognition

Check Readers

OCR is used to process banking checks without human involvement. A check can be inserted into a machine, the writing on it is scanned instantly, and the correct amount of money is transferred. A check reader captures check images and recognizes courtesy amounts and account information on the checks. Some readers also recognize the legal amount on the check and use information in both fields to cross-check the recognition results.

Bill Processing Systems

In general a bill processing system is used to read payment slips, utility bills and inventory documents.

The system focuses on certain region on a document where the expected information is located.

Airline Ticket Readers

In order to claim revenue from airline passenger ticket, an airline needs to have three records matched: reservation record, the travel agent record and the passenger ticket. However it is impossible to match all three records for every ticket sold. Several airlines are using a passenger revenue accounting system to account accurately for passenger revenue. The system reads the ticket number on a passenger ticket and matches it with the one in the airline reservation database. It scans tickets up to 260,000 tickets per day and achieves sorting rate of 17 tickets per second.

Handwritten Address Interpretation/ Address

Readers

Handwritten address interpretation technology has also been developed for address recognition systems for the United Kingdom's Royal Mail and the Australia Post. The address in a postal mail sorter locates the destination address block on a mail piece and read the ZIP code in the address book. If additional fields in the address block are read with high confidence the system may generate a 9 digit ZIP code is used to generate a bar code which is sprayed on the envelope.

General Purpose Page Readers

They are designed to handle a broad range of documents such as business letters, technical writing and newspapers. These systems capture an image of a document page and separate the page into text and non text region. Non-text regions such as graphics and line drawings are often saved separately from the text and associated recognition results. Text regions are segmented into lines, words and characters and the characters are passed to the recognizer. Recognition results are output in a format that can be post-processed by the application software. Most of these page readers can read machine written text only; a few can read hand-printed alphanumeric (Khaled *et al.*, 2015).

Limitations of Character Recognition System

OCR is unable to achieve a 100% recognition rate. Because of this, a system which permits fast and accurate recognition is a major requirement. The success of any OCR device to read accurately is the responsibility of the hardware manufacturer as well as depends on the quality of the items to be processed. (Kumar and Bhatiya, 2013) The main purpose of OCR from many years is as follows:

- To increase the accuracy of recognizing.
- To eliminate the need for specially designed fonts (character).
- In fact OCR is a time saver as it reduces the manual work, but it is not perfect.
- It hardly reaches more than 99.9% accuracy level.
- It faces problem with early printed books, newspaper etc.
- It faces problems with heavily bound material

Genetic Algorithms

Genetic algorithms are a family of computational models inspired by evolution. These algorithms encode a potential solution to a specific problem on a single chromosome and apply recombination operators to them so as to preserve critical information. GAs are often viewed as function optimizers, although the range of problems to which GAs have been applied is quite broad. The major reason for GAs popularity in various search and optimization problems is its global perspective, wide spread applicability and inherent parallelism. GA starts with a number of solutions known as population. These solutions are represented using a string coding of fixed length. After evaluating each chromosome using a fitness function and assigning a fitness value, three different operators- selection, crossover and mutation- are applied to update the population. An iteration of these three operators is known as a generation. If a termination criterion is not satisfied this process repeats. This termination criterion can be defined as reaching a predefined time limit or number of generations or population convergence (Kumar and Bhatiya, 2013).

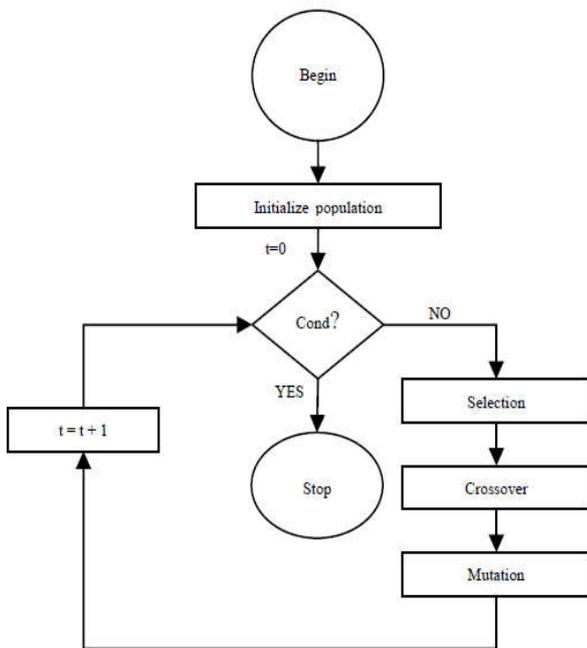


Fig. 2. A flowchart of working principles of a simple GA [7]

As it can be seen in Fig. 2 selection is the first operator applied on a population and forms a mating pool. Crossover operator is applied next to the strings of mating pool. It picks two strings from the pool at random and exchanges some portion of the strings between them. Mutation operator changes a 1 to 0 and vice versa. There are a number of factors that affect GA performance including appropriate operators, fitness function and population size. (Oliveira *et al.*, 2001) Inspired by the processes in biological evolution, it is based on recombination, natural selection, inheritance, recombination and mutation. Random populations are generated for which the fitness function is calculated in order to find the fitness function that is the most optimal. If the first set does not contain a fitness function value that is not satisfied, then chromosomes recombine among themselves and mutate to find a collection of another random population and the process continues.

The random population generated after recombination and selection is called the offspring. Genetic algorithm is basically a search technique that is faster than the classical ways of searching. They belong to evolutionary class of algorithm also known as EA. Each individual in a population of a genetic algorithm consists of properties (or traits) called the chromosome. These chromosomes can be mutated and altered. Chromosomes are generally represented in 0s and 1s. The population generation is an iterative process. Genetic operators (mutation and crossovers) are used to generate the next generation. Crossover is the process of combining the traits from two or more chromosomes in order to create a new chromosome. One point, two point are the most common crossover techniques. Other than these, techniques such as cut and splice and uniform crossover and half uniform crossover are also present. For selecting the chromosomes for crossover, again methods such as the Boltzmann selection or the tournament selection are use.

There are also a lot of other methods but it is beyond the scope of this paper to go into the very detail of each and every selection method. To show for an e.g. a crossover operation, consider the chromosomes 111111 and 111111. Chromosomes are crossed over to create another chromosome say 111101. This is an example of a single point crossover method. Mutation is to change the bits of the chromosomes to create a new chromosome. For e.g. chromosome strings can also be mutated, say 111101 to 011111 by changing just two bits and it gives a new chromosome. Different mutation types are available such as the bit string mutation, flip bit, boundary, non-uniform, uniform and Gaussian. (Yang and Honavar, 1998)

The algorithm terminates when:

- A solution has been found
- Time/resources drain out
- A fixed number of generation have been calculated
- There is no possibility of finding a better fitness function value.
- manually

The basic main advantages of using genetic algorithms are:

- Does not get stuck in a local optimum solution since it takes into account a population and not just one solution.
- It can be transferred to existing simulation and models.
- A large number of parameters are used and hence an accurate result can be found.
- The solution is not of fixed length.

Genetic algorithm based classification and recognition techniques

Vedgupt Saraf and Rao in (Vedgupt Saraf and Rao, 2013) have used genetic algorithm in character recognition of devnagari script. A flow chart for their entire process can be founded. Through their work using genetic algorithm, they claim to have an accuracy of around 97%-98%, although there are pairs that they found confusing. Pier Luca Lanzim *et al.* (Pier Luca Lanzim *et al.*, 1997), have used genetic algorithm for fast feature selection. The main advantage of their approach was that lesser processing time of CPU was required and the method was independent from a specific learning algorithm.

An example and the calculation method of the inconsistency rate can be found in (Yang and Honavar, 1998). They claim that using this method the algorithm is at least 10 times faster than a general genetic algorithm based feature selection. Ved Prakash Agnihotri (in Ved Prakash Agnihotri, 2010) presents the use of genetic algorithm in character recognition of handwritten devanagari script. In the feature selection phase the image is divided into zones, and in total 54 features are extracted for use in the recognition phase. The zones algorithm used for feature extraction phase. He in his work has achieved an 85.78% match and 13.35% mismatch. Chomtip Pornpanomchai, Verachad Wongsawangtham, Sathean Jeungudomporn and Nannaphat Chatsumpun in their Paper (ChomtipPornpanomchai *et al.*, 2011) presents the use of genetic algorithm in character recognition of handwritten Thai characters. The training data set consisted of 8160 characters and testing data set consists of 840 characters. They achieved an accuracy of 88.17%. There was a 10.10% mismatch and 1.66% rejection. They achieved a speed of 0.4192 seconds per character. The experiment was also conducted for English characters but a poor percentage approximately 1.06% match was obtained.

Vellingiriraj and Balasubramanie *et al.* (Vellingiriraj and Balasubramanie, 2014) have done a similar work to the thai recognition. The change is only in the script. They have used ancient Tamil handwritten characters instead of thai characters. They have also implemented with the same strategy and the same fitness value function. Difficulties have been faced when cursive writing has been used. Shashank Mathur (in Shashank Mathur, 2013) have also used genetic algorithm to implement character recognition. The process that he has used is to take the input string and to take a string of the current population. If error is above 10%, then the crossover stage is repeated. Else if the error is nearby 5%, the chromosome is mutated. These way new chromosomes are generated and through their work, they have obtained an overall efficiency of 79.16%.

Conclusion

The paper discusses as to what an offline character recognition system is. This paper discusses the various steps used in the process giving an overview of the details of the steps. Along with the overview, references to the detailed description of the process have been given below. Focus of this paper is mainly in the classification phase. The classification stage is explained in an overview again with the focus mainly on the genetic algorithm used for this process. In the genetic algorithm section, an overview of the genetic algorithm along with the details of the algorithm followed by its advantages has been discussed. The various methods of genetic algorithm that used in different papers have been explained along with the results and the efficiencies that they have achieved through their work. The work of genetic algorithms in the field of character recognition classification stage with an immense advancement in the Thai and Devanagari script can be clearly noticed. Also the process has been used in other languages as well. The paper describes in detail the use of genetic operators and how they are used. There has been a massive advancement in this field and more work is still under research. The main advantage of using genetic algorithm is that it operates on a set of solutions rather than a single solution at a time. This prevents the algorithm from getting stuck in local minima.

REFERENCES

- Chomtip Pornpanomchai, Verachad Wongsawangtham, Satheanpong Jeungudomporn and Nannaphat Chatsumpun, 2011. "Thai Handwritten Recognition by genetic algorithm (THCRGA)", *IACSIT*, Vol.3, No.2, April 2011.
- Deb, K. 1998. "Genetic Algorithm in Search and Optimization: the Technique and Applications", Proc. International Workshop on Soft Computing and Intelligent Systems, pp. 58-87, Calcutta, India, 1998.
- Khaled M.G. Noaman, Jamil Abdulhameed M. Saif and Ibrahim A.A. Alqubati, 2015. "Optical Character Recognition Based on Genetic Algorithms", *Journal of Emerging Trends in Computing and Information Sciences*, Vol. 6, No. 4, ISSN 2079-8407, April 2015.
- Kumar, G. and Bhatiya, P.K. 2013. "Analytical Review of Preprocessing Techniques for Offline Handwritten Character Recognition", *International Journal of Advances in Engineering Sciences*, Vol.3 (3), July, e-ISSN: 2231-0347 Print-ISSN: 2231-2013.
- Oliveira, L. S., Benahmed, N., Sabourin, R., Bortolozzi, F. and Suen, C. Y. 2001. "Feature Subset Selection Using Genetic Algorithms for Handwritten Digit Recognition" Proc. XIV Brazilian Symposium on Computer Graphics and Image Processing (SIBGRAPI'01), P.362, 2001.
- Pal, U. and Chaudhari, B. B. "Indian Script Character Recognition: a Survey", *IJSER*, Vol. 37, Issue 9, pp. 1887-1899, 2004.
- Pal, U., Sharma, N., Wakabayashi, T. and Kimura, F. "Off-Line Handwritten Character Recognition of Devnagari Script", In Proc. 9th ICDAR, pp.496-500, 2007.
- Pal, U., Wakabayashi, T. and Kimura, F. 2009. "Comparative study of Devanagari handwritten character recognition using different features and classifiers", in Proc. 10th Conf. Document Anal. Recognit, pp. 1111-1115.
- Pier Luca Lanzi and Politecnico di Milano, 1997. "Fast feature selection with genetic algorithm: A filter approach", *IEEE*, 0-7803-3949-5/97, 1997.
- Shashank Mathur, 2003. "self-evolving character recognition using genetic operators", Internaternational Book Series, Information Science & Computing, 2003.
- Soryani, M. and Rafat, N. 2008. "Application of Genetic Algorithms to FeatureSubset Selection in a Farsi OCR", World Academy of Science, *Engineering and Technology*, Vol: 2 2008-06-29,2008.
- Vedgupt Saraf, D.S. and Rao, 2013. "Devnagari script character recognition using genetic algorithm for better efficiency", *IJSCE*, ISSN: 2231-2307, Volume-2, Issue-4, April 2013.
- VedPrakash Agnihotri, 2010. "offline handwritten Devanagiri script recognition", *IJSER*, Vol-5, Issue-3, June 2010.
- Vellingiriraj, E. K. and Balasubramanie, P. 2014. "Recognition of ancient Tamil handwritten characters in palm manuscripts using genetic algorithm", *IJCST*, ISSN: 0976-8491, Vol.5, SPL-1, Jan-March 2014.
- Yang, J. and Honavar, V. 1998. "Feature Subset Selection Using a Genetic Algorithm", Proc. IEEE Intelligent Systems, vol. 13, no. 2, pp. 44- 49, 1998.