

**A Review on Text Detection and Character Recognition for Printed Document**Aayushi J Solanki,¹, Haresh Chande²¹Computer Engineering Department, HJD Institute of Technical Education and Research, Kera-Kutch,²Computer Engineering Department, HJD Institute of Technical Education and Research, Kera-Kutch,

Abstract—Text recognition in images is a research area which attempts to develop a computer system with the ability to automatically read the text from images. These days there is a huge demand in storing the information available in paper documents in a format to a computer storage disk and then later reusing this information by searching process. One simple way to store information from these paper documents in a computer system is to first scan the documents and then store them as images. But to reuse this information it is very difficult to read the individual contents and searching the contents from these documents line-by-line and word-by-word. The challenges involved in this are the font characteristics of the characters in paper documents and quality of images. Due to these challenges, a computer is unable to recognize the characters while reading them. Thus there is a need of character recognition mechanisms to perform Document Image Analysis which transforms documents in paper format to electronic format. In this paper we have discussed method for text recognition from images.

Keywords- pre-processing, segmentation, character recognition, classification

I. INTRODUCTION

Text Recognition is the mechanical or electronic conversion of images of typed, handwritten or printed text into machine encoded text. It is widely used as a form of data entry from printed paper data records, whether passport documents, invoices, bank statements, computerized receipts, business cards, mail, printouts of static data, or any suitable documentation. It is a common method of digitizing printed texts so that it can be electronically edited, searched, stored more compactly, displayed on-line, and used in machine processes such as machine translation, text-to-speech, key data and text mining. OCR is a field of research in pattern recognition, artificial intelligence and computer vision.

Text recognition has become one of the most successful applications of technology in the field of pattern recognition and artificial intelligence. Many commercial systems for performing OCR exist for a variety of applications, although the machines are still not able to compete with human reading capabilities. Optical Character Recognition deals with the problem of recognizing optically processed characters. Now-a-days, there is growing demand for the software systems to recognize characters in a computer system when information is scanned through paper documents as we know that we have a number of newspapers and books which are in printed format related to different subjects. These days there is a huge demand in "storing the information available in these paper documents into a computer storage disk and then later reusing this information by searching process". One simple way to store information in these paper documents into a computer system is to first scan the documents. Whenever we scan the documents through the scanner, the documents are stored as images in a format in the computer system. These images containing text cannot be edited by the user.

The reason for this difficulty is the font characteristics of the characters in paper documents are different from the font of the characters in a computer system. As a result, a computer is unable to recognize the characters while reading them. This concept of storing the contents of paper documents in a computer storage place and then reading and searching the content is called document processing. Sometimes in this document processing we need to process the information that is related to languages other than the English in the world. This process is also called Document Image Analysis (DIA). Thus our need is to develop some text recognition algorithm to perform Document Image Analysis which transforms documents in paper format to electronic format.

II. LITERATURE SURVEY

A Neural Network Based Handwritten Meitei Mayek Alphabet Optical Character Recognition System [1]

Handwritten character recognition is a part of optical character (OCR) system. OCR can be applied to both printed text and handwritten documents. In this paper we discussed the handwritten character recognition of Meitei Mayek (Manipuri script). Although OCR has been studied and developed for many Indian scripts very few works have been reported so far for Meitei Mayek. This paper describes the handwritten Meitei Mayek (Manipuri script) alphabets recognition (HMMAR) using a neural network approach. The alphabet database is pre-processed and the extracted features are sent to a neural network system for training. The trained neural network is further tested and performance analysis is observed.

The emphasis is given on the process of character segmentation from a whole document i.e. isolating a single character image from a complete scanned document.

Text Recognition from Images [2]

Text recognition in images is a research area which attempts to develop a computer system with the ability to automatically read the text from images. These days there is a huge demand in storing the information available in paper documents in a format to a computer storage disk and then later reusing this information by searching process. One simple way to store information from these paper documents in a computer system is to first scan the documents and then store them as images. But to reuse this information it is very difficult to read the individual contents and searching the contents from these documents line-by-line and word-by-word. The challenges involved in this are the font characteristics of the characters in paper documents and quality of images. Due to these challenges, a computer is unable to recognize the characters while reading them. Thus there is a need of character recognition mechanisms to perform Document Image Analysis (DIA) which transforms documents in paper format to electronic format. In this paper we have discussed method for text recognition from images.

Recognition of Off-line Hand printed English Characters, Numerals and Special Symbols [3]

Optical Character Recognition can improve the interaction between man and machine in various applications including data entry, office automation, digital library, banking applications, health insurance and tax forms etc. Much of work has been done in the recognition of machine printed characters in various languages with considerably good efficiencies, however making robust recognition engines that can be put to recognize handwritten and hand printed data with commendable recognition rates still remains as an active area of research owing to the challenges like diverse human handwriting style, variation in shape, angle and style of characters. Taking into account the challenges and scope for improvement in this domain, the work of off-line character recognition of hand printed document images containing English Characters-Uppercase and Lowercase, Numerals and Special Characters has been presented. Statistical, Geometric and Directional Feature Extraction techniques have been applied over segmented character image. Classification was done using Multilayer perception neural network (NN) with back propagation and Support vector machine (SVM) classifier.

Neural Network based Handwritten Character Recognition system without feature extraction [4]

Handwriting recognition has been one of the active and challenging research areas in the field of image processing and pattern recognition. It has numerous applications which include, reading aid for blind, bank cheques and conversion of any handwritten document into structural text form. In this paper an attempt is made to recognize handwritten characters for English alphabets without feature extraction using multilayer Feed Forward neural network. Each character dataset contains 26 alphabets. Fifty different character datasets are used for training the neural network. The trained network is used for classification and recognition.

Scanning Neural Network for Text Line Recognition [5]

This paper describes segmentation free text line recognition approach using multilayer perceptron (MLP) and hidden Markov models (HMMs). A line scanning neural network trained with character level contextual information and a special garbage class is used to extract class probabilities at every pixel succession. The output of this scanning neural network is decoded by HMMs to provide character level recognition.

A Robust Approach for Offline English Character Recognition [6]

In this paper, we propose a recognition model based on Artificial Neural Network (ANN) supported by novel feature extraction technique. Handwritten data has continued to persist as a means of recording information in day-to-day life with the introduction of latest technologies. The constant development of computer tools lead to the requirement of easier interface between human and computers. Recognition of handwritten characters by computer is a complicated task as compared to typed character. The proposed system is implemented using MATLAB successfully. The ANN accepts the input as a scanned image. This input undergoes a sequence of pre-processing steps; binarization and normalization. Then features are extracted and matched from the stored data in the database.

Optical Character Recognition [7]

Character recognition techniques associate a symbolic identity with the image of a character. In a typical OCR system input characters are digitized by an optical scanner. Each character is then located and segmented, and the resulting character image is fed into a pre-processor for noise reduction and normalization. Certain characteristics are then extracted from the character for classification.

III. GENERAL OVERVIEW

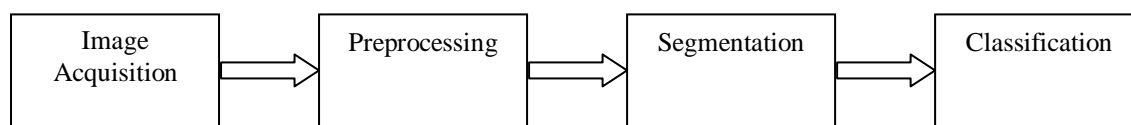


Figure.1- A general overview

Character Recognition proceed in mainly four steps.

(1) Image Acquisition (2) Preprocessing (3) Character Segmentation and (4) Classification.

IV. IMAGE ACQUISITION

In Image acquisition, the recognition system acquires a scanned image as an input image. The images should have a specific format such as JPEG, BMP etc. This image is acquired through a scanner, digital camera or any other suitable digital input device.

A. Image Pre-processing

The process of enhancing the image, which should be used for further processing, is called preprocessing. Preprocessing is the major step in handwriting recognition system. Noise in a document image is due to poorly photocopied pages. These scanned images not only have noises which are inbuilt within it, but also the noise may be during the scanning of that image.

a) Gray scale processing

It includes converting RGB to gray for noise removal, border enhancement for brightness. Colored images are also complex in space and time, therefore it is necessary to convert them to gray scale to reduce time and space complexity. The basic idea behind gray conversion is to eliminate hue and saturation by not affecting its luminance. For this, we compute the threshold of an image by using the suitable gray scale value. This separates the object of interest from background. Thresholding is important to provide sufficient contrast of an image such that, varying level of intensity between foreground and background are considered. Gray scale conversion enhances the quality of an image for later computational processes. Gray scale images consist of different ranges of gray values; from 0 to 255.

b) Binarization

The conversion of the gray scale image to black and white is called binarization. Binary images are also called as Bi-level or two-level. First, the original RGB images should be converted to gray scale and then the image is converted to black and white image. Most of the OCR packages work on the binarized images. The conversion is possible because of the threshold values and the values which are higher than the threshold are white and the values which are lower than this threshold are black. Otsu's method is used to perform threshold based on cluster i.e. from gray level image to binary image.

c) Normalization

The process of changing the intensity value of the pixel to the range of [0, 1] is called normalization in image processing. The conversion of various dimension images into fixed dimensions is also called as normalization.

d) Thinning

Thinning is a pre-process which results in single pixel width image to recognize the character easily. It is applied repeatedly leaving only pixel-wide linear representations of the image characters.

e) Median Filtering

It comes under the category of nonlinear filters. It changes the gray value of the pixels to the median of the gray value of surrounding pixels. We use a 3x3 mask and calculate the corresponding gray value of each pixel using the 8 neighboring pixels. This helps in noise removal. Median filtering gives advantages such as no reduction in contrast since output values are its neighborhood values, boundaries remain unchanged. Median filters are very useful in the presence of impulse noises also called salt and pepper noise because of its appearance as white and black dots superimposed on an image.

f) Morphological Operations

They are generally used to remove noise from the imperfect segmentation. Morphological operations are especially suited for binary images. So they are performed on output image of thresholding. Here Dilation and erosion are performed. Dilation and erosion are used to remove holes in the detected foreground. In the process of dilation the size and shape determination of structuring element is very important

V. SEGMENTATION

A) Line Segmentation

This segmentation is the most important process in text recognition. Segmentation is done to make the separation between the individual characters of an image. Line segmentation process is segmenting the text region into lines, also called as line segmentation.

B) Character Segmentation

Word segmentation is easier than line segmentation and character segmentation. Space between two words is generally more than three pixels. Words are segmented by the projection based method.

C) Skew Detection

The skew detection used in character recognition where in the whole word is inscribed in a polygon with at least two dimensions. The center of gravity of the image is considered as a single line extending at a certain angle with the horizontal. The angle is measured which gives the angle by which the word or document is rotated and also signifies the direction and angle by which it should be rotated for it to be a text in readable and normal form.

D) Hough Transform

This is a technique for feature detection. Initially it was used for line detection. Now it has been extended to find position of different shapes like circle, square, oval etc. Therefore we can use it to detect number plates from the image based on their rectangular shape.

E) Bounding Box

Segmentation is one of the most important processes in the number plate recognition, because all further steps rely on it. If the segmentation fails, a character can be improperly divided into two pieces, or two characters. The ultimate solution on this problem is to use bounding box technique. The bounding box is used to measure the properties of the image region. Once a bounding box created over each character and numbers presented on image, each character & number is separate out for recognition.

VI. CLASSIFICATION

The classification is the process of identifying each character and assigning it to the correct character class, so that texts in images are converted into computer understandable form.

A) Template matching

It is useful when we have characters of fixed size. This technique is different from the others in that no features are actually extracted. Instead the matrix containing the image of the input character is directly matched with a set of prototype characters representing each possible class. The distance between the pattern and each prototype is computed, and the class of the prototype giving the best match is assigned to the pattern.

B) Neural network

Character Recognition algorithm can also be used to recognize the vehicle number. The character can be recognized using neural network like artificial neural network (ANN), back propagation neural network, kohonen neural network etc. In the Artificial Neural Network (ANN) for classification neural network can get itself trained automatically on the basis of efficient tools for learning large databases and examples. This approach is non-algorithmic and trainable. The Kohonen neural network works differently than the feed forward neural network. The Kohonen neural network contains only an input and output layer of neurons. There is no hidden layer in a Kohonen neural network.

C) Support Vector Machine (SVM)

SVMs have become more and more important in the field of pattern recognition. SVM is forcefully competing with many methods for classification. An SVM is a supervised learning technique. SVM takes Statistical Learning Theory (SLT) as its theoretical foundation, and the structural risk minimization as its optimal object to realize the best generalization. They are based on some simple ideas and provide a clear intuition of what learning from examples is all about. More importantly, they possess the feature of high performance in practical applications. The SVMs use hyper planes to separate the different classes. Many hyper planes are fitted to separate the classes, but there is only one optimal separating hyper plane. The optimal one is expected to generalize well in comparison to the others. A new data sample is classified by the SVM according to the decision boundary defined by the hyper plane. Among many classification methods, SVM has demonstrated superior performance.

VII. CONCLUSION

Character recognition has different phases and accuracy of each phase dependent on previous phase. Different kinds of page or worst document and skew lines or characters under which the systems are designed to operate are the main challenges. From this survey it can be concluded that various pre-processing methods like median filter, morphological operation that are erosion and dilation used

to improve the image quality. Segmentation can be done by various techniques like line and word segmentation. For recognition process, various classification methods like neural network, SVM can be used to improve the result.

REFERENCES

- [1] Romesh Laishram, Pheiroijam Bebison Singh, Thokchom Suka Deba Singh and Sapam Anilkumar, Angom Umakanta Singh, "A Neural Network Based Handwritten Meitei Mayek Alphabet Optical Character Recognition System", 2014 IEEE International Conference on Computational Intelligence and Computing Research.
- [2] Mr. Pratik Madhukar Manwatka, Mr. Shashank H. Yadav Department, "Text Recognition from Images", *IEEE Sponsored 2nd International Conference on Innovations in Information, Embedded and Communications Systems (ICIIECS) 2015*
- [3] Nisha Sharma, Bhupendra Kumar, Vandita Singh, "Recognition of Off-line Hand printed English Characters, "Numerals and Special Symbols", *2014 5th International Conference-Confluence The Next Generation Information Technology Summit (Confluence)*
- [4] J. Pradeep, E. Srinivasan, S. Himavathi, "Neural Network based Handwritten Character Recognition system without feature extraction", *International Conference on Computer, Communication and Electrical Technology – ICCET 2011, 18th & 19th March, 2011*
- [5] Sheikh Faisal Rashi, Faisal Shafait and Thomas M. Breuel, "Scanning Neural Network for Text Line Recognition", *10th IAPR International Workshop on Document Analysis Systems 2012*
- [6] Suman Avdhesh Yadav, Smita Sharma, Shipra Ravi Kumar, "A Robust Approach For Offline English Character Recognition, 2015 1st International conference on futuristic trend in computational analysis and knowledge management (ABLAZE 2015)
- [7] Jagruti Chandarana¹, Mayank Kapadia, "Optical Character Recognition", *International Journal of Emerging Technology and Advanced Engineering (ISSN 2250-2459, ISO 9001:2008 Certified Journal, Volume 4, Issue 5, May 2014)*