

Text Recognition in an Image Using Statistical Method

Dr. Manmohan Singh
Associate Professor (CSE&IT)
RKDF (SOE), Indore

Ms. Megha Singh
Assistant Professor (CSE&IT)
CIIT, Indore

Mr. Rahul Sharma
Assistant Professor (CSE&IT)
RKDF (SOE), Indore

Abstract— Nowadays content-based retrieval is pivotal issue in design of multimedia systems, which requires pattern matching and human cognition process. Literature survey reveals that huge scope is available to work in the area of text information retrieval from image files. Due to varying pattern in the image size, shape and orientations of text in images, it is computationally difficult to extract information from image files. The work is based on the Image Processing techniques and cognitive graphics. In this work an image will be pre-processed through the image processing techniques, and later will be matched with the standard available structures. This work will be a small attempt to understand the working of algorithms to identify and extract the individual letters and numbers.

Keywords— *Multimedia, Pattern Matching, Information Retrieval, Cognitive Processes, Image Processing*

I. INTRODUCTION

The Human Brain, One of the most powerful processor of the world has the incredible capability to identify, recognize understand and to differentiate. The vision or say perception of brain in its own is a unique feature of nature Our brain would easily identify recognize and differentiate between different types of shape size, orientation, textures, colors and so it could easily observe the unique and different patterns. However, when it comes to computational ability a machine takes over the advantage. The text in an image could be seen and recognized by the brain easily but if we say it for a computer it becomes a tough task. The computer should be able to recognize the characters in an image by some way or the other, which may be of any type, so the text, could be retrieved and thus opens a way to perform word processing over the text.

II. REVIEW OF CONCEPTS AND THEORIES

Following are some state of art techniques in Information retrieval from multimedia systems,

A. Classification and Learning for Character Recognition:

Classification methods is based on learning from examples. These techniques have been widely applied to character recognition from the 1990s. It includes statistical methods, artificial neural networks, support vector machines, multiple classifier combination etc

B. Optical Character Recognition:

This technique was developed to translate scanned images of handwritten, typewritten or printed text into machine- encoded text. Tesseract, originally developed as proprietary software at Hewlett-Packard between 1985 and 1995, now sponsored by Google, is considered to be one of the most accurate open source OCR engine currently available.[2]

C. Text Correction:

Since the result returned by the OCR engine is not be always correct due to image imperfections. Hence The text correction is a necessary step after OCR .This type of errors can be categorized into so called non-word error – which means that the text string returned by OCR does not correspond to any valid word in a given word set. Existing robust text correction algorithms have a good performance in correcting this type of non-word error.[2]

D. Text Extraction:

Text extraction techniques are widely studied because text embedded in images and videos may provide important information. Many characteristics of text regions have been summarized and characterized effectively by several features, e.g. text pixels have near-homogeneous colour, character strokes form distinct texture, etc. Y. Hasan and J. Karam developed a text extraction algorithm that utilized morphological edge/gradient detection. [2]

E. Artificial neural networks:

Pattern recognition is the area that employees' techniques like Feed forward neural networks, including multilayer perceptron (MLP), radial basis function (RBF) network, higher-order neural network (HONN) etc. Here the connecting weights are usually adjusted to minimize the squared error on training samples in supervised learning. Using a modular network for each class was shown to improve the classification accuracy. A network using local connection and shared weights, called convolution neural network, has reported great success in character recognition [6].

III. PROPOSED METHOD

From Literature survey it is clear that less significant work has been done in the field of text extraction from a multimedia content. Some work has been done in identifying text contents from image files. We propose here a method for recognizing each and every character as a different object. At the initial level the method will count the number of characters present in the Image and individually check them with the predefined dataset.

A. Assumptions: Consider an image with following characteristics

1. RGB image containing only two colors, black and white.(For example scanned copy of a B&W newspaper)
2. Proposed algorithm will Works only for special font type.
3. In the current paper we have have taken font type in image as Arial, Arial Black.
4. Proposed algorithm works for font size 10 to 72 .

The proposed system has the following key points:

1. A rectangular boundary or window that surrounds characters in the image file is selected. this reduces the processing overheads.
2. A mechanism to count the number of objects (characters) in sample RGB image is to be devised.
3. Each character is individually checked against a fixed value from predefined dataset.
4. co-relation coefficient among the datasets with the input text part of RGB image will be used to check the matching percentage. The higher the correlation coefficient shows higher matching accuracy that leads to high probability of the specified characters.

The above proposed method which recognize characters can be divided in following steps

- Segmentation Algorithm
- Character Detection algorithm
- Binarizing Technique
- Matching Algorithm

B. Detection Algorithm:

The detection algorithm detects the local text area, by first reducing an RGB image X to a gray scale image Y in relation with $Y = 0.299R + 0.587G + 0.114B$ Where each pixel value is calculated by above formula. The gray scale image has colors (shades of gray) ranging between 0-255. A thresh holding value from 100-210 be taken to reduce the higher color values. Later the connected components are found and labelled accordingly, after that the text area including all the connected component only is returned. The text area is selected by finding the height of connected components closer to the boundary. This algorithm also counts the number of connected components.

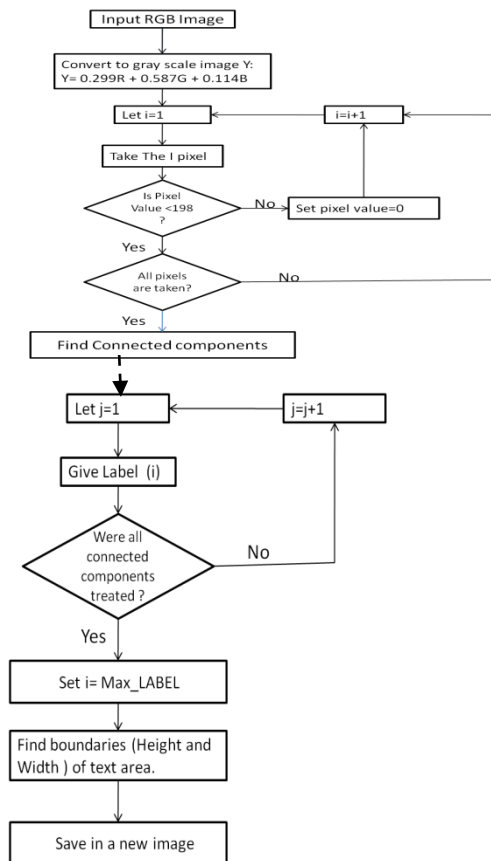


Fig. 1- Detection algorithm

C. Binarization:

In this technique the image is represented by only two values 0 and 1. The 0 represents the absence of the data whereas the 1 represents the presence of data. Binarization is an important step serving the purpose of only remaining text(data) in an image.

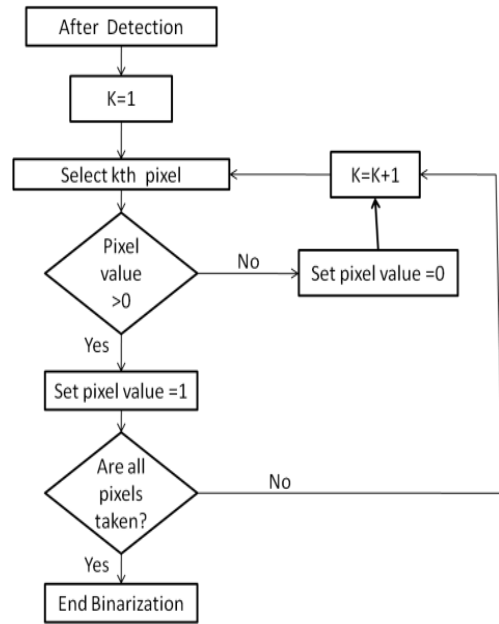


Fig 2 Binarization Algorithm

D. Segmentation:

The segmentation allows each of the objects in the text area to be cut and saved in a structure for further processing. Each object can now be seen in seen with its real characteristics.

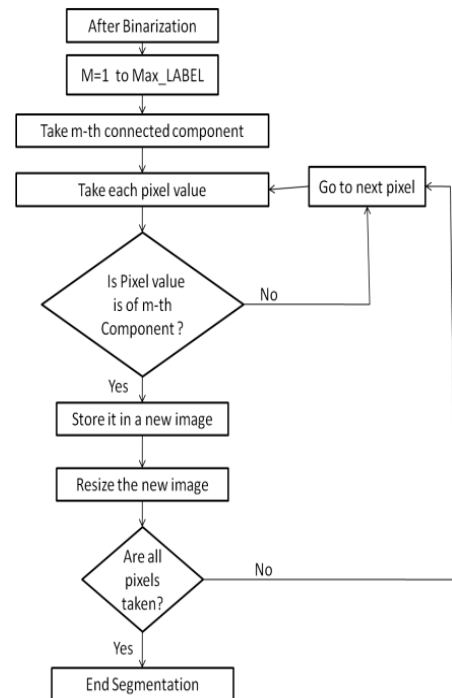


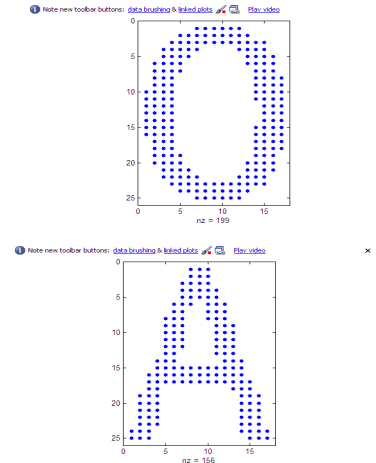
Fig 3 . Segmentation

E. Matching Algorithm:

The resulting template is matched with standard templates and on the basis of correlation coefficient for the two $N \times M$ matrix (25 x 17 in above case) type of object structures each having pixel values as there elements

V. ANALYSIS OF DATA

The data we have collected is the data having the image representation as bits representation and by observing the bit pattern the data can be identified. Some of the graphical patterns like spy graph, contours, and normal bit representation of the image were plotted of which spy graphs are shown below for some random digits and alphabet like 0 and A. The other characters are not represented here but have the same kind of bit representations.



CONCLUSIONS

The work includes 4 steps: Detection, Binarizing, segmentation followed by matching. Applying a simplistic statistical method of correlation provided an advantage of solving certain critical identification issues such as detection of character 'I', the difference between alphabet 'C', 'Q', 'O', and the most important one difference between alphabet 'O' and digit '0'. The accuracy of the current work is about 99 % (approx.) for font face Arial ,however size may vary from 10 to 72.the RGB image taken in account included two colors, black and white. This work had a large scope and opportunities in future in the field of content based retrieval, computer vision, machine learning, and artificial intelligence.

Acknowledgment

We would like to convey our sincere thanks to Dr. R.K. Vyas, Director IIPS Devi Ahilya University Indore for his out of bound support , motivation throughout the time and providing facilities at Development Center for our R&D activities. We would also like to our support amber jain , sachin saxena and ajeet khan for always being constant source of inspiration for us. It is only because of his vision the work came into life. We are also very thankful to him for his valuable guidance, suggestions and teaching in all fields.

References

- [1] Mohanad Alata — Mohammad Al-Shabi, Text detection and character recognition using fuzzy image processing, Journal of ELECTRICAL ENGINEERING, VOL. 57, NO. 5, 2006, 258–267
- [2] Derek Ma, Qiuhan Lin, Mobile camera based text detection and translation Department of Electrical Engineering and Tong Zhang from Department of Mechanical Engineering from Stanford University
- [3] Datong Chen, Jean- Marc Odobez, Herv/e , Text detection and recognition in images and video frames Boulard from Dalle molle Institute for Perceptual Artificial Intelligence (IDIAP), Pattern Recognition 37 (2004) 595 – 608.
- [4] Fakhreddin Mamedov, Jamal Fathi Abu Hasna,Character

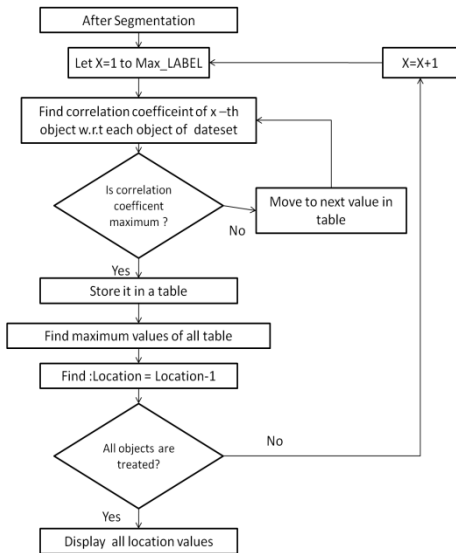
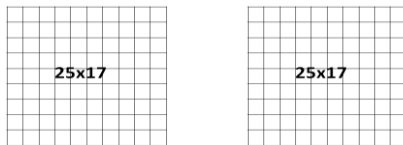


Fig 4. Matching Algorithm.

The method followed in the in the proposed system was based on calculating correlation coefficient for the two matrix with same dimension where the elements of matrix are the corresponding pixel value.



Matrix 1 (Predicted Structure) Matrix 2 (Standard Structure)

In order to estimate the character in found structure the matrix of type 1 should be compared with each standard available structure.

IV. COLLECTION OF DATA

Data for 0:

0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0
0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0
0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0
0	0	1	1	1	1	1	1	0	0	0	0	0	0	1	1	1	1	0	0	0
0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	1	1	1	1	0	0
0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	1	1	1	1	0
0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0
0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0

Data for 'A':

0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	0
1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

recognition using neural networks.

- [5] Dimitri Van De Ville ,Noise reduction by fuzzy image filtering , Member, IEEE, Mike Nachtegaele, Dietrich Van der Weken, Etienne E. Kerre, Wilfried Philips, Member, IEEE, and Ignace Lemahieu, Senior Member, IEEE, IEEE Transaction on Fuzzy Systems, Vol. 11,No.. 4, August 2003
- [6] Cheng-Lin Liu , Classification and learning for character recognition: Comparison of Methods and Remaining Problems by National Laboratory of Pattern Recognition (NLPR) Institute of Automation, Chinese Academy of Sciences and Hiromichi Fujisawa Central Research Laboratory, Hitachi, Ltd.