

Offline Handwritten & Typewritten Character Recognition using Template Matching

Sunny Kumar*

School of Computer Science & Engineering
Bahra University
Solon, India
csunny56@yahoo.com

Pratibha Sharma*

School of Computer Science & Engineering
Bahra University
Solon, India
pratibhasharma80@gmail.com

Abstract: In the modern times identity verification has become an important task for various reasons. Mostly identity provided by the authorities is the combination of characters & digits. Now there are several techniques available to recognize them efficiently, & good accuracy rates also have been achieved for characters & digits recognition (Typewritten). But in case of handwritten characters, accuracy rate becomes low. In this paper performance of template matching technique is analyzed for handwritten & typewritten character recognition using two parameters i.e. accuracy rate & time taken for execution. Here template matching technique has been applied to offline Typewritten character set & offline Handwritten characters set.

Keywords-OCRS, segmentation, binarisation, classification, matching, templates

I. INTRODUCTION

Character recognition is a process of converting handwritten, typewritten or printed text images into machine encoded code or text [3]. Character recognition is the research field of pattern recognition, artificial intelligence & computer vision [1]. Character recognition has a variety of applications such as data entry for business documents, automatic number plate recognition & make electronic images of the printed documents such as Google books ([1],[5],[4],[6]).

A. Online Character Recognition

Online character recognition is the real time recognition of characters. Since Online character recognition uses online systems which have better timing information for recognizing characters. Online Character Recognition also avoids the initial step of locating the characters ([1],[4],[5]).

B. Offline Character Recognition

Offline character recognition involves the automatic conversion of text from an image into letter code. In this type of character recognition, Handwritten/Typewritten characters usually scanned from images & then converted into gray/binary scale image & then fed to recognition algorithm. Offline Character recognition is more challenging task than online. Since in this type of recognition we have no control over the medium & devices ([1],[4],[5]).

II. OFFLINE CHARACTER RECOGNITION SYSTEM

The major steps of offline character recognition system used in this project are:-

A. Gray scaling

In this phase, input image is converted into gray scale image basically gray scale image is a black & white image in which each pixel is considered a single sample which carries intensity information. In gray scale digital image shades of gray varying from black at weakest intensity & white at strongest ([7],[2]).

B. Binarisation

In this phase gray scale image is converted into a binary image. In this paper a binary method called Otsu algorithm is used for binarisation. Binarisation is done Block by Block[8].

1. Otsu method

In Otsu's method a threshold is needed that minimizes the intra-class variance (the variance within the class), defined as a weighted sum of variances of the two classes:

$$\sigma_w^2(t) = \omega_1(t)\sigma_1^2(t) + \omega_2(t)\sigma_2^2(t) \quad (1)$$

Where ω_i denotes the probabilities of the two classes separated by a threshold t , and σ_i denotes the variances of these classes.

$$\sigma_b^2(t) = \sigma^2 - \sigma_w^2(t) = \omega_1(t)\omega_2(t) [\mu_1(t) - \mu_2(t)]^2 \quad (2)$$

Which is expressed in terms of class probabilities ω_i and class means μ_i , which in turn can be updated iteratively.

2. Otsu Algorithm

i) Compute histogram and probabilities of each intensity level.

ii) Set up initial $\omega_i(0)$ and $\mu_i(0)$

iii) Step through all possible thresholds $t = 1 \dots$ maximum intensity.

1. Update ω_i and μ_i

2. Compute $\sigma_b^2(t)$

iii) Desired threshold corresponds to the maximum $\sigma_b^2(t)$.

iv) You can compute two maxima (and two corresponding thresholds). $\sigma_{b1}^2(t)$ is the greater max and $\sigma_{b2}^2(t)$ is the greater or equal maximum.

v) Now image binarisation can be expressed as

$$I_B(x, y) = \{0 | I(x, y) < t^* / 1 | I(x, y) \geq t^* \quad (3)$$

Here 0 denotes black representing the text, and 1 denotes white representing the background [8].

C. Filtering

The median filtering is a nonlinear digital filtering technique often used to remove the noise. Image filter can be linear or nonlinear. Here two dimensional median filter is used to remove the dust or noise from the image. Two dimensional median filtering is implemented on an image using a mask of odd length, the mask moves over the image and at each center pixel the median value of the data within the window is taken as the output [9].

D. Cropping

Now image is cropped fit to text. After that each line separates from rest of lines.

E. Connected Component Labeling

Connected component labeling can be applied on binary or gray level images and different measures of connectivity are possible. Connectivity checks are carried out by checking the labels of pixels that are North-East, North, North-West and West of the current pixel. (assuming 8-connectivity).

A Faster algorithm for connected component extraction [10]

On the first pass:

1. Iterate through each element of the data by column, then by row (Raster Scanning)
2. If the element is not the background
 - i) Get the neighboring elements of the current element
 - ii) If there are no neighbors, uniquely label the current element and continue
 - iii) Otherwise, find the neighbor with the smallest label and assign it to the current element
3. Store the equivalence between neighboring labels

On the second pass:

1. Iterate through each element of the data by column, then by row
2. If the element is not the background
 - i) Relabel the element with the lowest equivalent label

Here, the background is a classification, specific to the data, used to distinguish salient elements from the foreground. If the background variable is omitted, then the two-pass algorithm will treat the background as another region [10].

Example of two-pass algorithm

8-connectivity based [11]

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

Figure 1: Digital Image from which connected components are to be extracted

After the first pass.

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	1	0	0	2	2	0	0	3	3	0	0	4	4	0
0	1	1	1	1	1	1	1	1	0	0	3	3	3	3	0	0
0	0	0	1	1	1	1	0	0	0	3	3	3	3	0	0	0
0	0	1	1	1	1	0	0	0	3	3	3	0	0	3	3	0
0	1	1	1	0	0	1	1	0	0	0	3	3	3	0	0	0
0	0	1	1	0	0	0	0	0	5	3	0	0	0	3	3	0
0	0	0	0	0	0	6	6	5	3	0	0	7	3	3	3	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 2: Digital Image with Different labels

A total of 7 labels are generated

The label equivalence relationships generated are

Set ID	Equivalent Labels
1	1,2
2	1,2
3	3,4,5,6,7
4	3,4,5,6,7
5	3,4,5,6,7
6	3,4,5,6,7
7	3,4,5,6,7

Figure 3: Set IDs of generated labels

Labels generated after the merging of labels is carried out. Here, the label value that was the smallest for a given region "floods" throughout the connected region and gives two distinct labels, and hence two distinct labels

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	1	0	0	1	1	0	0	3	3	0	0	3	3	0
0	1	1	1	1	1	1	1	1	0	0	3	3	3	3	0	0
0	0	0	1	1	1	1	0	0	0	3	3	3	3	0	0	0
0	0	1	1	1	1	0	0	0	3	3	3	0	0	3	3	0
0	1	1	1	0	0	1	1	0	0	0	3	3	3	0	0	0
0	0	1	1	0	0	0	0	0	3	3	0	0	0	3	3	0
0	0	0	0	0	0	3	3	3	3	0	0	3	3	3	3	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 4: Digital image with two labels

F. Segmentation

Image segmentation is a process of splitting the digital image into multiple segments. The main aim of the segmentation is to simplify the digital image or change the representation of an image for some useful meaning. Image segmentation is used to locate the boundaries & objects in images[12].

G. Classification

Now normalization process is applied on input images. Each point will be normalized into 42 x 24 pixels this is the standard size to perform the 2d-correlation. Correlation in two dimensions is the main operation in classification process. Correlation is used to determine the likeness of the point of entry to the workforce. It will give the similarity value between two matrices or images.

III. TEMPLATE MATCHING

Template Matching is a technique used for mapping one image into another. Template matching has a variety of applications such as character recognition, object recognition, and classification etc. In this project we have designed 7 samples of handwritten & typewritten templates. Size of each image is fixed i.e. 24*42. BMP image format has been used for designing the templates having a bit depth of 1 bit. These samples are firstly loaded into MATLAB individually, and then all handwritten samples are combined to make a single sample. Same procedure is used for creating single sample of typewritten templates. By doing this we have seven samples of each character & each digit for matching.

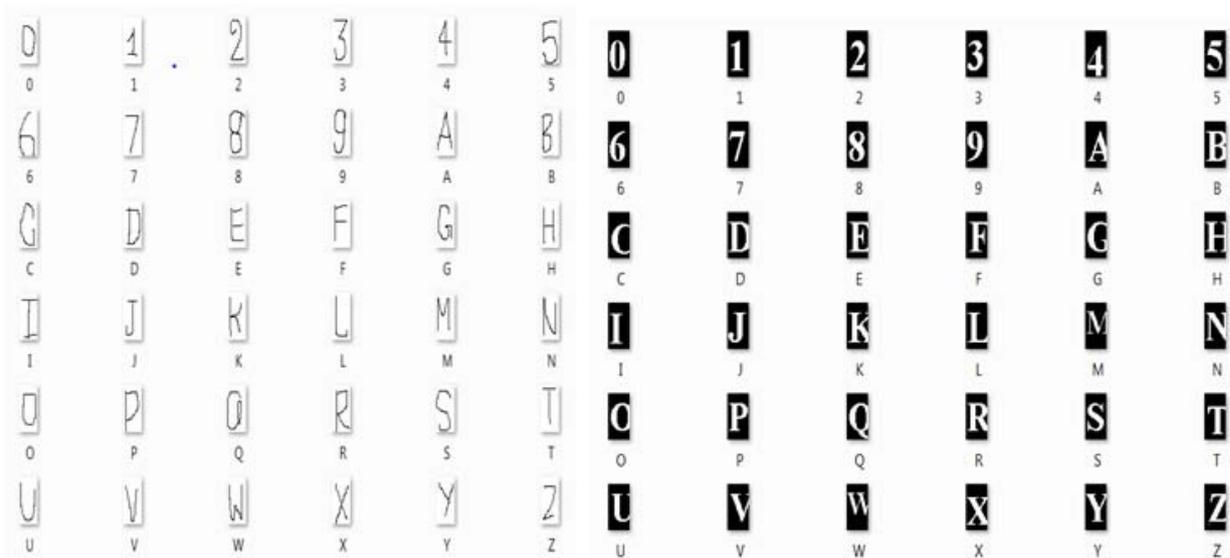


Fig 5:- A Sample of handwritten & typewritten templates from seven samples

IV. IMPLEMENTATION AND RESULTS

This project is implemented in Matlab R2009a & then 360 images of characters & digits are tested (180 for Handwritten & 180 for Typewritten). Five images of each character is taken for testing purpose & each image contains an average of 4 characters. Accuracy rate & time taken for simulation will be the two parameters on which template matching technique is analyzed for handwritten & typewritten characters. Accuracy is calculated by no. of characters recognized to the total number of characters in an image.

For example, if we have following information

Parameter	Image1	Image2	Image3	Image4	Image5	Average
Accuracy	0.80	1	0.60	0.50	0.34	0.648
Time taken	0.39	0.34	0.36	0.36	0.35	0.36

Then accuracy for an image is = No. of characters recognized/ total no. of characters.

$$0.80+1+0.60+0.50+0.34/5 = 0.648 \text{ (accuracy rate for a character)}$$

$$0.39+0.34+0.36+0.36+0.35/5 = 0.36 \text{ (average time taken by a image containing 5 characters)}$$

Now average time taken by each character is $0.36/5 = 0.07$

In this way accuracy rate & average time taken for each character is calculated for 26 upper case letters & 10 numbers for different cases. These cases are given below.

1. Typewritten Templates & Typewritten Characters: - In this case typewritten characters are matched with typewritten templates.
2. Typewritten Templates & Handwritten Characters: - handwritten characters are matched with typewritten templates
3. Handwritten Templates & Typewritten Characters: - typewritten characters are matched with handwritten templates
4. Handwritten Templates & Handwritten Characters:-handwritten characters are matched with handwritten templates.

These results for these cases are shown in figures

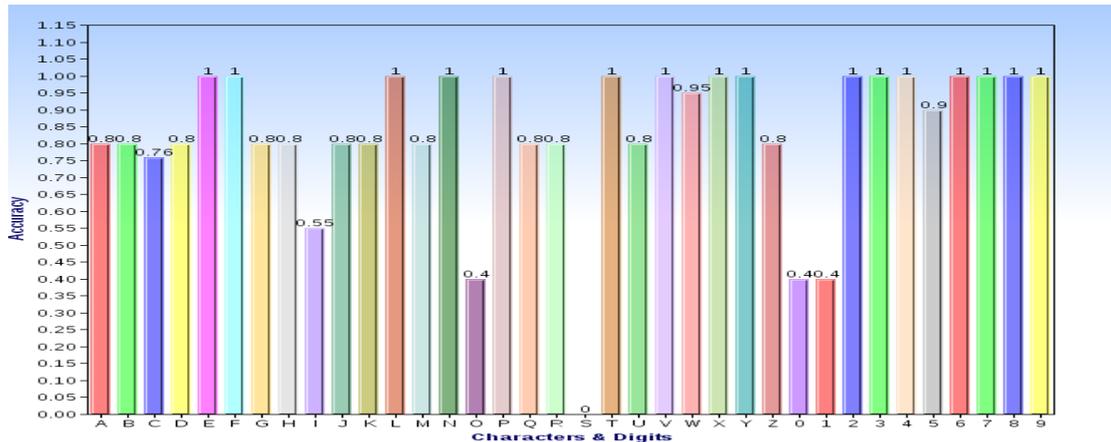


Fig 6: Accuracy rate for typewritten templates & typewritten characters

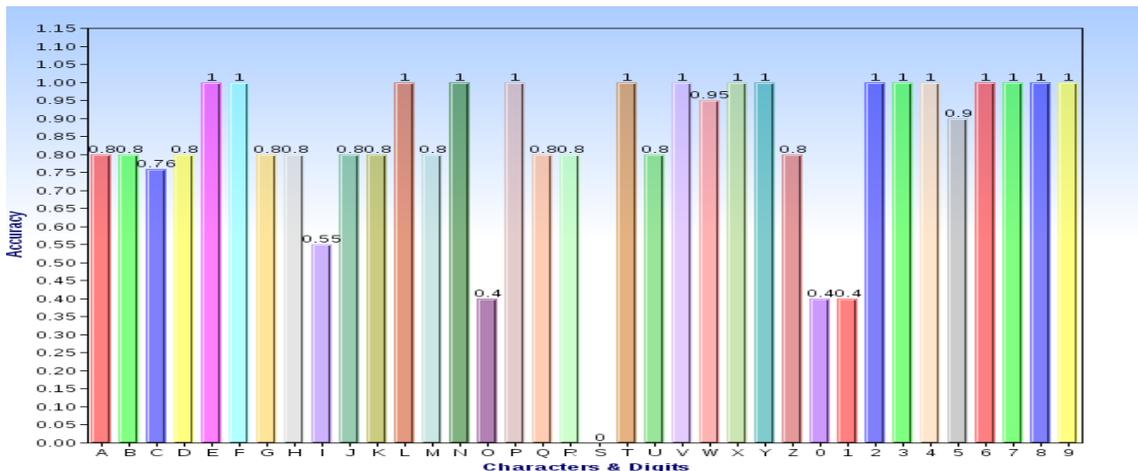


Fig 7: Accuracy rate for typewritten templates & typewritten characters

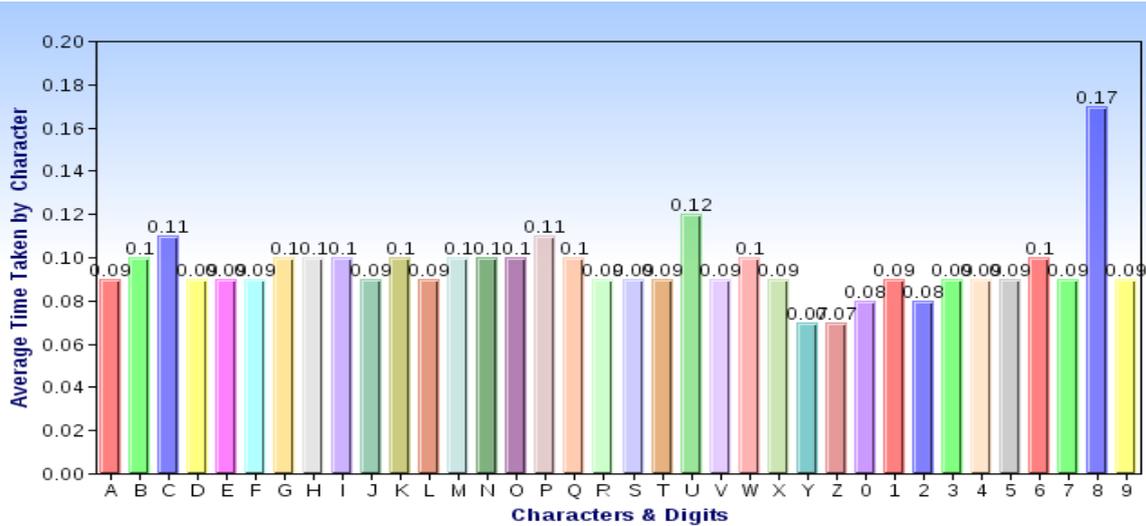


Fig 8: Average time taken by character for typewritten templates & typewritten characters

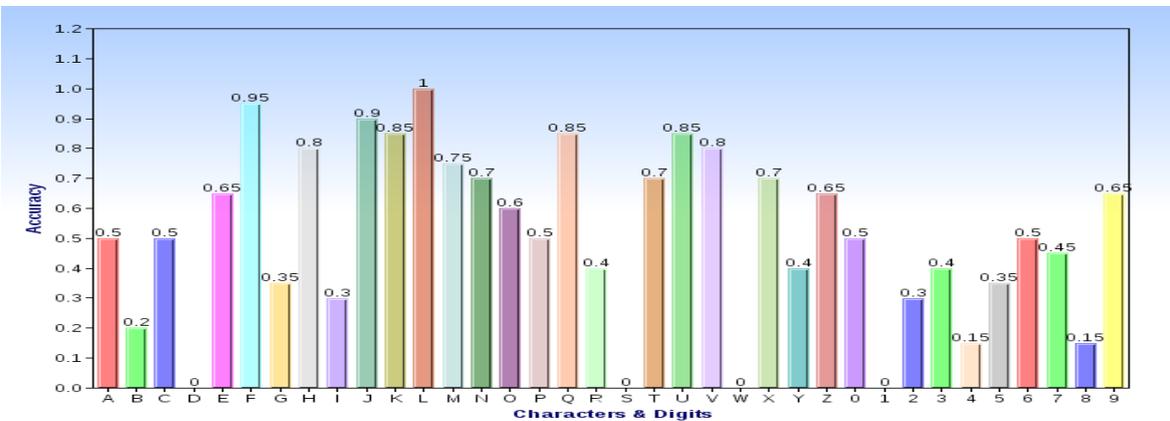


Fig 9: Accuracy rate for typewritten templates & handwritten characters

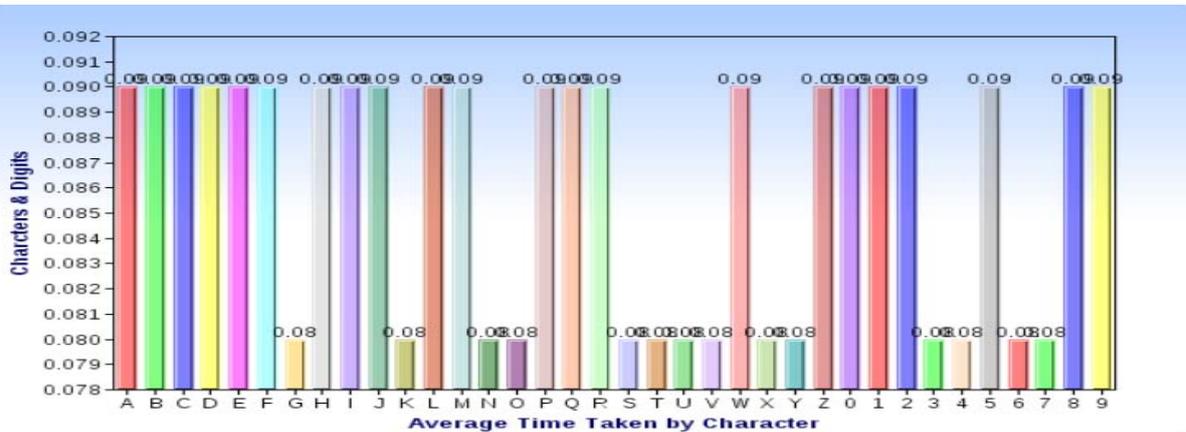


Fig 10: Average time by character for typewritten templates & handwritten characters

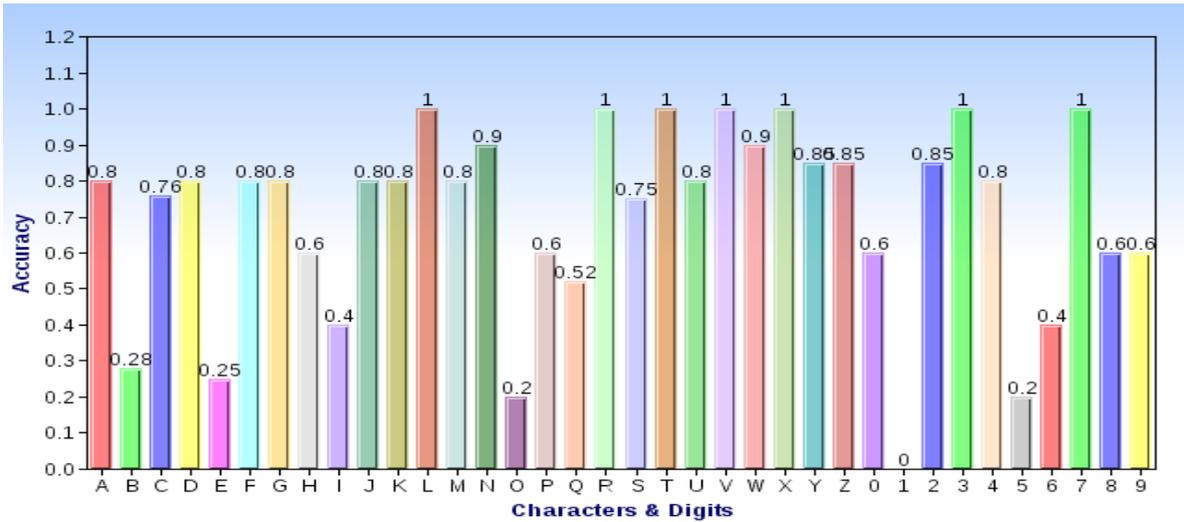


Fig 11: Accuracy rate for handwritten templates & typewritten characters

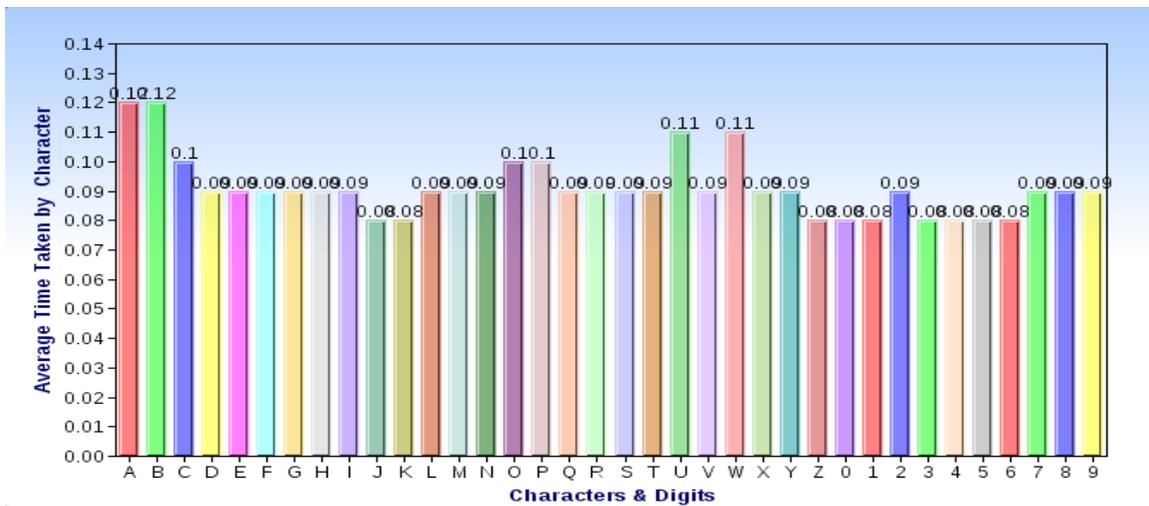


Fig 12: Average time taken by character for handwritten templates & typewritten characters

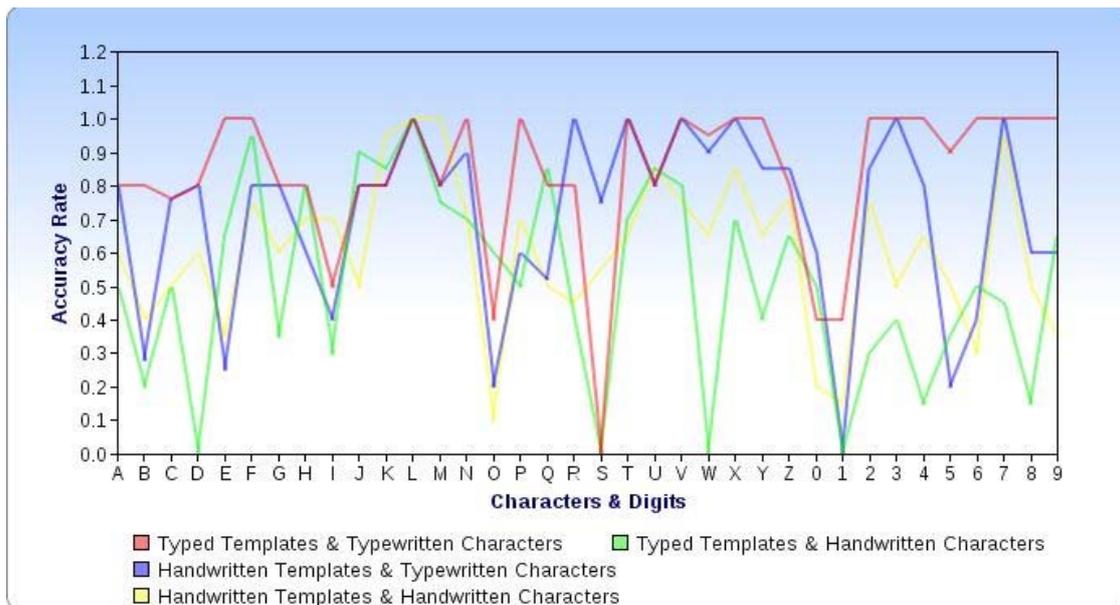


Figure 13: Combined graph showing accuracy

From the above graph it has been observed that typewritten templates & typewritten characters has higher accuracy rate than all others & typewritten templates & handwritten characters has lowest accuracy rate

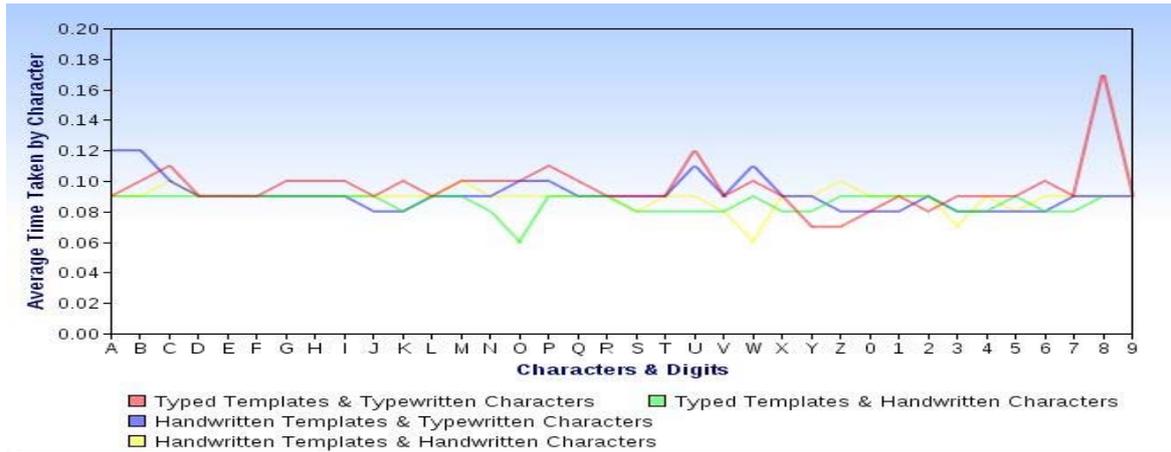


Figure 14: Combined graph showing time taken by character

Now overall accuracy of a particular set is calculated by combining all the accuracy rate obtained for characters & digits and then that value is divided by 36 to find the average accuracy rate for a system. In the same way time taken for characters is calculated. On the basis of results obtained following accuracy rates have been achieved. For typewritten templates & typewritten characters accuracy rate is 83%. For typewritten templates & handwritten characters accuracy rate is 50%. For handwritten templates & typewritten characters accuracy rate is 70%. For handwritten templates & handwritten characters accuracy rate is 60%. Time taken for simulation is approximately same for all the cases. On the basis of accuracy rates obtained things are saying that typewritten templates for typewritten characters has better accuracy rate than handwritten templates for handwritten characters & opposite section means typewritten templates for handwritten characters & handwritten templates for typewritten characters giving accuracy ranges from 50-60%. So, opposite section is not giving good results as expected.

V. CONCLUSION

In this paper we tried to compare the performance of Template matching technique. The performance of the Template Matching was evaluated in terms of Recognition accuracy rate & Time taken for execution. The experimental results obtained from the study shows that the template matching is better technique for Typed Character recognition than Handwritten Character recognition. But if well shaped templates are designed using some electronic media such as electric notepad, pen computing device & input image containing good efficient handwriting style. Then Accuracy rate for handwritten characters can be increased using template matching.

VI. REFERENCES

- [1] Priya Sharma and Randhir Singh "Performance of English Character Recognition with and without Noise" International Journal of Computer Trends and Technology- volume 4 Issue 3- 2013
- [2] Debasish Biswas, Amitava Nag, Anjan Pal, Soumadip Ghosh, Arindrajit Pal, Sushanta Biswas and Snehashish Banerjee "Novel gray scale conversion Techniques based on Pixel depth" Journal of Global Research in Computer Science, Volume 2, No. 6, June 2011
- [3] Schantz, Herbert F. (1982). The history of OCR, optical character recognition. [Manchester Center, Vt.]: Recognition Technologies Users Association. ISBN 9780943072012
- [4] http://www.cs.uic.edu/~srizvi/BIT_Thesis.pdf
- [5] Special Issue on Character Recognition and Document Understanding, IEICE Trans. Information and Systems, vol. E79-D, no. 5, July 1996
- [6] Special Issue on Oriental Character Recognition, Pattern Recognition, vol. 30, no. 8, 1997
- [7] Stephen Johnson. Stephen Johnson on Digital Photography. O'Reilly. ISBN 0-596-52370-X
- [8] Divya gilly, Dr. Kumudha raimond / International Journal of Engineering Research and Applications (IJERA) ISSN: 2248-9622 www.ijera.com Vol. 3, Issue 2, March -April 2013, pp.1240-1245
- [9] V. MUSOKO AND A. PROCHAZKA "NON-LINEAR MEDIAN FILTERING OF BIOMEDICAL IMAGES" INSTITUTE OF CHEMICAL TECHNOLOGY, DEPARTMENT OF COMPUTING AND CONTROL ENGINEERING
- [10] Lifeng He; Yuyan Chao; ; Suzuki, K. (1 May 2008). "A Run-Based Two-Scan Labeling Algorithm". IEEE Transactions on Image Processing 17 (5): 749-756. doi:10.1109/TIP.2008.919369. PMID 18390379
- [11] R. Fisher, S. Perkins, A. Walker and E. Wolfart (2003). "Connected Component Labeling
- [12] http://elearning.vtu.ac.in/17/e-Notes/DIP/segmentation_DIP-SDG.pdf
- [13] M. Cheriet, "Extraction of Handwritten Data From Noisy Gray-Level Images Using A Multiscale Approach", Pattern Recognition and Artificial Intelligence, Vol 13, No. 5 (1999), pp. 665-684
- [14] http://portal.fke.utm.my/fklibrary/files/faraharinarosli/2012/362_FARAHARINAROSLI2012.pdf