# Enhanced Handwritten Document Recognition Using Confusion Matrix Analysis

Umadevi T P [a,1] and Murugan A [b]

[a] *Assistant Professor, Department of Computer Science, JBAS College for Women (Autonomous), Chennai, India*
[b] *Associate Professor & Head, PG & Research, Department of Computer Science, Dr. Ambedkar Government Arts College (Autonomous) Affiliated to University of Madras, Chennai, India*

**Abstract.** The handwritten Multilanguage phase is the preprocessing phase that improves the image quality for better identification in the system. The main goals of preprocessing are diodes, noise suppression and line cancellation. After word processing, various attribute extraction techniques are used to process attribute properties for the identification process. Smoothing plays an important role in character recognition. The partitioning process in the word distribution strategy can be divided into global and local texts. The writer does not use this header line to write the text which creates a problem for skew correction, classification and recognition. The dataset used are HWSC and TST1. The tensor flow method is used to estimate the consistency of confusion matrix for the enhancement of the text recognition .The accuracy of the proposed method is 98%.

**Keywords.** Handwritten, Morphological, tensor flow, Optical Character Recognizer, angular directions.

## 1. Introduction

Automatic document processing helps conversion of paper document into digital text form. Many OCR algorithms has been developed by the researchers for different scripts. But the performance of OCR algorithm becomes unsatisfactory. The challenges associated with the processing of handwritten [1,2] document image are due to the irregular handwriting.

When two characters touch each other, a large space set is formed between the numbers .This is very important for separation, as the points of contact are often close. First, the size and shape of the text is analyzed and detected when touched. Finally, depending on the touching position of the text and the morphology of the contact area with the incision site are new algorithm also formed [3]. It has been observed that the identification of manuscripts with structured functions based on preprocessing documents is effi-

---

[1]Umadevi TP, Department of Computer Science, JBAS College for Women (Autonomous), Chennai.
E-Mail: Umashiva06@gmail.com

cient. Therefore, word-level structural features [4,5] are used for text identification. First, the words in the document image are sorted into text boxes [6,7]. In the final process, sometimes two areas of text form the words. These text correction rules are divided into words. Coherence analysis is performed to allow the possibility of general conditions in the pixel pairs of the altered image [8,9]. Based on the calculated design characteristics, visuals were defined using the K-Nearest Neighbor (K-NN) classification. These design features define the script without knowing its composition. The advantage lies in the design features, not the design features.
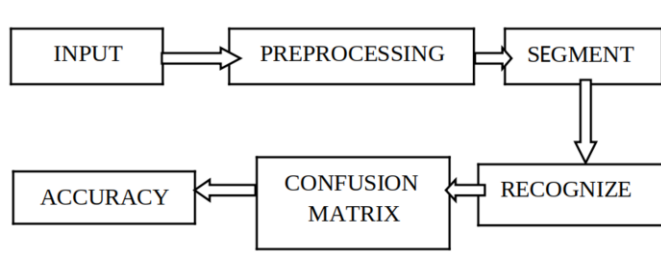


**Figure 1.  Flow of work**

## 2.  Literature Review

Wenyu Zhang et al. [10] introduced a new anti-mission learning (AML) paradigm to improve HCR performance. Data in limited training, includes prior knowledge of printed data and is independent of the author features of semantics. Of the available manual methods, the AML method offers a different one authors use independent semantic functions automatically as knowledge prior to standard print data does realistic research. To address issues of speed and memory capacity, Xiao et al. [11,12] introduced Globalization Controlled range extension method and optimal weight loss Tensor flow method . The proposed method is evaluated in a database dedicated to a publicly available database .

## 3.  The Proposed Recognition System

### 3.1.  Preprocessing

This step includes document digitization, followed by binarization. We have used a HP flat bed scanner for digitizing the collected documents in 300 dpi. Initially all the images are stored at 256 intensity level or as grey level. During binarization procedure the grey (256 level) images are converted to binary (2 level) image or binary image. Preprocessing helps us to improve the character recognition system [13]. Preprocessing is essential in order to have a higher recognition rate. There is certain constraints in hand written document. The handwriting should be legible and uniform. No decorative or cursive style letter should be written. Letters should remain specific beyond every mean of system. Thickness over the line about the slip must be specific.

## 3.2. Algorithm Classification For Handwritten

This research work have used different steps in the classification process belief-level fission techniques is used in which each classifier use new methods of tensor flow .This is unlikely to be a style course, a confidence score of 1 indicates that the test sample format is the maximum perhaps the corresponding script is a trust score results of the individual classifiers are collected and validated [14]. The highest scoring script is marked as an internal category. In our experiment, the research work collected results from SVM, KNN (K = 5), and neural network classification [15,16] for SVM classification [17], Belief is built on personal level works. It is calculated as a percentage of losses. The output value of a node belonging to a particular class that assigns a belief value to the nervous system. Classified build-in classification can achieve accuracy 98% of 5x cross valid data is standard deviation. Algorithm

## 4. Implementation

The character recognition system of a handwritten includes steps such as digitization, pre-validation process and classification. Handwriting recognition system is widely adopted. It depends on the functions to be removed. The removed attributes must be able to classify the character habit [18,19]. In this application, the research work have provided a great advantage by using a horizontal cover there is also a maximum limit for functions. This research work also compared the results of evaluating the proposed system with new systems. Electronic translation is done using a process that contains symbols image is scanned and is an electronic representation of the real image of the character in TIFF format. The image has been created. Digitization creates a digital image that is provided during image preprocessing. The grayscale image is simplified in a 360 x 360 window using the nearest neighbor Interpolation Algorithm (IA). After simplifying, the research work creates a bitmap for the entire image. Now, the bitmap was converted into a thinner image using the parallel attenuation algorithm proposed by Zhang. To determine the style displayed in each image, the research work used the vertical, stroke and shadow functions features of functions and basic functions.
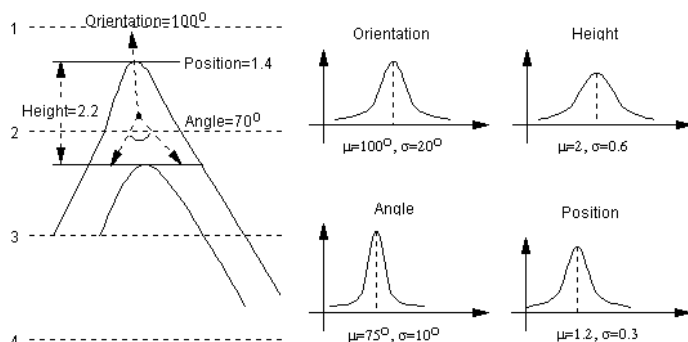


**Figure 2.  Angle of Character**

---

**Algorithm 1 Classification Handwritten Algorithm**

---

1: **procedure** IMAGE DOCUMENT(Model S, Dataset P)
2:      Sheet = create a knot ()
3:      $pixel.test_c ondition = findbestsplit(s, f)$
4:      $V = v|Possibleresultof vroot.test_c ondition$
5:      VεP for any value=0 to 360 degree
6:      HW+ Image document = Image document (HW, P)
7:      Handwritten. current = head
8:      **while** *current ≠ null* **do**
9:          current = current. Next
10:          new Dataset= new Handwritten.(o)
11:      N.next = head→head = new.dataset;
12:      **if** tail == null) **then**
13:          tail = head;
14:          head=new Handwritten.(element)
15:      **else if** size == 1 **then**
16:          Handwritten. temp = head
17:      Handwritten. current = head;
18:      **return** temp.element (Filter)
19: **procedure** REMOVE(index)
20:      **if** *index < 0||index >= size* **then**
21:          **return** null
22:      **else if** index == 0 **then**
23:          **return** remove First(HCAR)
24:      **else if** index == size - 1 **then**
25:          **return** remove Last(character)
26:      Handwritten. previous = head;
27:      **for** int i = 1; i < index; i++ **do**
28:          previous = previous. Next(TST1);
29:          Handwritten. current = previous. Next;
30:          previous. Next = current. Next;
31:      **if** element value < the value in current. Element **then**
32:          parent = current. Left
33:      **return** true(value)

---

## 5. Result and Evaluation

To train tattoo network we used the datasets of HWSC and TST1 databases. The number of eligible classes is 3755.

$$P(V;W) = \sum_H \frac{1}{Z} exp - E(V, H; W)$$
$$E(V, H; W) = \sum_k E(V, H_k; W_k)$$
$$E(V, H_k; W_k) = -H_k \otimes Filter(V, W_k)$$

This research work is implemented our scheme on 12,000 numbers collected by different people from different of schools, colleges, universities, teachers, banks and post employees, traders etc. The data set includes the writing style. Note that the accuracy

of the definition scheme is 92.8%.It is not enough to know works number. There are two types of businesses based on this nervous system. Numerical identification using the proposed identification scheme are represented by dotted lines. (X1, y1) and (x2, y2) are integers the first and third rolls of the dataset collection in the dataset have 98% accuracy MLP based project proposed by Bhattacharya et al. accuracy is 39.83% was tested using 3,330 points. Therefore, our proposed scheme gives better results than tensor flow based classifiers. During the experiment, the research work found that a high definition rate of (eight) is 97.8%. That is the reason, confusion rates are calculated depends on the post-thinning groups can be fraught with difficulties [8]. They were also in the queue due to the different styles of writing; the following method does not work as expected there is a moving number divided into the figure thinning and uniform morphology method ($\sigma$) .Ours the identification method accurately identifies the number of this part without any change System. The disadvantage of the proposed method is that it does not give results if there is a discrepancy in the small part used as an outline collection range. However, if these faulty points are properly connected, performance can be improved received. To eliminate some of these conditions when using the size of the stopping area.

$$P(H_k|V) = \sigma(Filter(V, H_k))$$
$$P(V|H) = \sigma(\sum_k Filter(H_k, W_k))$$
$$\sigma(x) = \frac{1}{1 + exp - x}$$

The data HCAR is in the standard Tier 1 kit and has hundreds of samples of each track. All test roles are in training Time, which is the identity of the closed group. However, there are about 7000 characters in the group organization for modern travel.

**Table 1.  Confusion matrix data analysis of handwritten predictions**

| Measure | Formula | Percentage |
|---|---|---|
| Sensitivity | TPR = $\frac{TP}{TP+FN}$ | 88% |
| Specificity | SPC = $\frac{TN}{FP+TN}$ | 82% |
| Positive Predictive Value (Precision) | PPV = $\frac{TP}{TP+FP}$ | 78% |
| Negative Predictive Value | NPV = $\frac{TN}{TN+FN}$ | 84% |
| False Positive Rate | FPR = $\frac{FP}{FP+TN}$ | 66% |
| False Discovery Rate | FDR = $\frac{FP}{FP+TP}$ | 56% |
| False Negative Rate | FNR = $\frac{FN}{FN+TP}$ | 44% |
| Accuracy | ACC = $\frac{TP+TN}{TP+TN+FP+FN}$ | 97.88% |
| F1 Score | F1 = $\frac{2TP}{2TP+FP+FN}$ | 91.22% |
| Matthews Correlation Coefficient | MCC = $\frac{TPxTN - FPxFN}{\sqrt{(TP+FP)x(TP+FN)x(TN+FP)x(TN+FN))}}$ | 89.1% |

## 6.  Conclusion

In this paper, an approach to distinguish machine printed and handwritten text from document image is proposed. Text separation helps better processing of real life documents by applying separate OCR on separate type of text images. Simple and easy to compute

feature vector is constructed for the same. The present work is focused on English, Chinese and Arabic. Experimentation is done on these three scripts individually script wise and in collectively combining all three scripts. The system requires a table processing program that provides a complete alphabet of printed or handwritten letters by scanning of new n algorithms using tensor flow given 98% better results.

## References

[1]   Gupta D, Bag S. CNN-based multilingual handwritten numeral recognition: A fusion-free approach. Expert Systems with Applications. 2021 Mar 1;165:113784.

[2]   Pal U, Chaudhuri BB. Automatic separation of machine-printed and hand-written text lines. InProceedings of the Fifth International Conference on Document Analysis and Recognition. ICDAR'99 (Cat. No. PR00318) 1999 Sep 22 (pp. 645-648). IEEE.

[3]   Kavallieratou E, Stamatatos S. Discrimination of machine-printed from handwritten text using simple structural characteristics. InProceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004. 2004 Aug 26 (Vol. 1, pp. 437-440). IEEE.

[4]   Jang SI, Jeong SH, Nam YS. Classification of machine-printed and handwritten addresses on korean mail piece images using geometric features. InProceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004. 2004 Aug 26 (Vol. 2, pp. 383-386). IEEE.

[5]   Zheng Y, Li H, Doermann D. Machine printed text and handwriting identification in noisy document images. IEEE transactions on pattern analysis and machine intelligence. 2004 Jun 28;26(3):337-53.

[6]   Raghuraman G, Sabena S, Sairamesh L. Image retrieval using relative location of multiple ROIS. Asian Journal of Information Technology. 2016;15(4):772-5.

[7]   Sharieff AH, Sabena S, Sathiyavathi V, SaiRamesh L. Intelligent framework for joint data hiding and compression using SMVQ and fast local image in-painting. Int. J. Sci Technol. Res. 2020;9(2):2267-71.

[8]   Imade S, Tatsuta S, Wada T. Segmentation and classification for mixed text/image documents using neural network. InProceedings of 2nd International Conference on Document Analysis and Recognition (ICDAR'93) 1993 Oct 20 (pp. 930-934). IEEE.

[9]   Kumar A, Awasthi N. An efficient algorithm for text localization and extraction in complex video text images. In2013 2nd International Conference on Information Management in the Knowledge Economy 2013 Dec 19 (pp. 14-19). IEEE.

[10]  Zhang W, Yang D, Zhang S, Ablanedo-Rosas JH, Wu X, Lou Y. A novel multi-stage ensemble model with enhanced outlier adaptation for credit scoring. Expert Systems with Applications. 2021 Mar 1;165:113-123.

[11]  Li Z, Wu Q, Xiao Y, Jin M, Lu H. Deep matching network for handwritten Chinese character recognition. Pattern Recognition. 2020 Nov 1;107:107-17.

[12]  Li Z, Wu Q, Xiao Y, Jin M, Lu H. Deep matching network for handwritten Chinese character recognition. Pattern Recognition. 2020 Nov 1;107:107471.

[13]  Kuhnke K, Simoncini L, Kovacs-V ZM. A system for machine-written and hand-written character distinction. InProceedings of 3rd International Conference on Document Analysis and Recognition 1995 Aug 14 (Vol. 2, pp. 811-814). IEEE.

[14]  Anand J, Sivachandar K. An edge vector and edge map based boundary detection in medical images. International Journal of Innovative Research in Computer and Communication Engineering. 2013 Jun;1(4):1050-1055.

[15]  Simard PY, Steinkraus D, Platt JC. Best practices for convolutional neural networks applied to visual document analysis. InIcdar 2003 Aug 6 (Vol. 3, No. 2003).

[16]  Kunte RS, Samuel RS. On-line character recognition for handwritten kannada characters using wavelet features and neural classifier. IETE Journal of Research. 2000 Sep 1;46(5):387-93.

[17]  Ramteke SP, Gurjar AA, Deshmukh DS. A novel weighted SVM classifier based on SCA for handwritten marathi character recognition. IETE Journal of Research. 2019 Jun 20:1-3.

[18]  Sahare P, Chaudhari RE, Dhok SB. Word level multi-script identification using curvelet transform in log-polar domain. IETE Journal of Research. 2019 May 4;65(3):410-32.

[19]  Cao Z, Lu J, Cui S, Zhang C. Zero-shot Handwritten Chinese Character Recognition with hierarchical decomposition embedding. Pattern Recognition. 2020 Nov 1;107:107488.