

Literature Survey on Writer Identification for Handwritten Document Based on Structured Learning

Kiran Gole¹, Julekha Mulani², Govind Kumar³, Vishal Sherekar⁴

^{1,2,3,4} ISB & M School of Technology Nande Village, Tal. Mulashi, Pune Savitribai Phule Pune University, India

Abstract

OCR (optical character recognition) is an essential task for word segmentation. The features of fonts present in handwritten document are irregular and those are different depending on the person then, it is considered a challenging problem. To overcome this problem, we formulating the problem of word segmentation as a binary quadratic assignment problem which considers pair wise correlations between the gaps as well as of individual gaps in the document. We are using the Structured SVM (Support Vector Machine) framework which will estimate all the parameter to work the proposed method well regardless of different writing styles and written languages without user-defined parameters. And using Word segmentation we will make the task of identifying the writer's of document easy.

Keywords: Handwritten Document, word segmentation, SVM.

1. Introduction

Segmentation of the image is important task to match the features of two different images. And hence to understand and also to match the image of handwritten document the segmentation of document into the word and text-line is the important task. Unlike the machine-printed document, the handwritten document is more challenging because of-

- (i) irregular spacing presents in the word
- (ii) variations of writing styles of person

To solve this problem, this paper proposes a scale invariant feature transform (SIFT) Algorithm which extracts the features of corresponding document's image and help to identify the writer of document. Identifying the writer is essential task in today's world

2. Literature survey

Language Independent Text-Line Extraction Algorithm for Handwritten Documents

This paper proposes a language-independent text-line extraction algorithm for the processing of handwritten document images. By introducing stroke lengths, we split under-segmented CCs into several pieces to have better representations for text components. Then with the help of text line extraction we can estimate the line spacing and orientation of every CC. A system that will trace out trained data from input file. Input file is consisting of handwritten text and printed text. Text is separated using CC segment and then next in text line extraction, that estimate the line spaces and overlapping text lines. And uses Optical Character Recognition algorithm. Resulting keywords are matched with the trained dataset. The keywords which pass the set threshold level of frequency will be added in output file as a result of handwritten recognition system. Trained dataset contains number of different samples of each handwritten character.

Handwritten Text Segmentation using Average Longest Path Algorithm.

In this paper a new global non- holistic method for handwritten text segmentation, this does not make any limiting assumptions on the character size and the number of characters in a word. Specifically, the proposed method finds the text segmentation with the maximum average likeliness for the resulting characters. For this purpose, we use a graph model that describes the possible locations for segmenting neighboring characters, and we then develop an average longest path algorithm to identify the globally optimal segmentation. Uniformly and densely sample the text image to construct a set of candidate segmentation boundaries. A directed graph is then constructed to embed the character likeliness of the text segments between any two candidate segmentation boundaries. In this graph, text segmentation is reduced to the problem of finding an average longest path between the first and the last candidate segmentation boundary. It find that the average longest path in the constructed graph can be found in polynomial time with global optimality.

A New Segmentation Algorithm for Online Handwritten Word Recognition in Persian Script.

In this paper a new segmentation algorithm for the main stroke of online Persian handwritten words. Using this segmentation, which present a perturbation method which is used to generate artificial samples from handwritten words. Our recognition system is composed of three modules. The first module deals with the preprocessing of the data. Also propose a wavelet-based smoothing technique which enhances the recognition performance compared to the conventional widely used technique. The second module is word segmentation into convex portions of the global shape which we call Convex Curve Sectors (CCSs). The third module is to analyze those CCSs and use the information for recognition performed by Dynamic Time Warping (DTW) technique. Using CCSs provides the DTW-based classifier with a compact word representation which makes comparison much faster.

3. Architecture Design

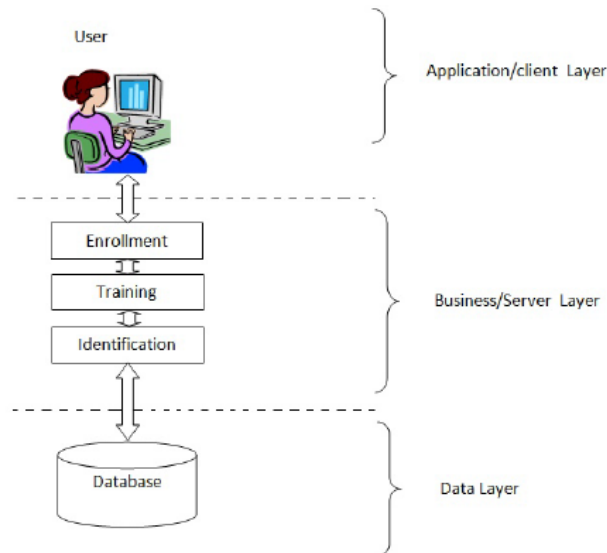
The problem will be solved using SIFT algorithm, by dividing the solution in three parts these are:

1. Enrollment:

The phase in which user will get register and receive the User name and Password.

2. Training:

The user will access the system using provided id and password, and give image as the input to calculate and extract the features of handwritten document .



3. Identify:

In this phase the all the algorithms will get applied on new image to extract the features of image ,and to match those features with existing one for indentifying the writer.

4. Algorithm description

The main contribute of the paper is SIFT algorithm, this algorithm used to extract the features of handwritten document's image.

For filtering the image some algorithms are applied

1. Ostu's Binary image :This Function is used to convert the image into the grayscale image and the grayscale image into binary image
2. Remove blur: This function filter out the image by removing the blur present in image.
3. Isotropic LOG(Laplacian Of Gaussian) filter: This function used for edge detection, it performs following steps on image
 1. Start processing with image.
 2. Blur the given image .
 3. Perform laplacian on image which is blurred.
 4. Find out zero crossing of laplacian then at this point to threshold compare the local variance ,declare edge if the threshold get exceeded.
 5. Calculate median of the filter image.

4. Structure Invariant Feature Transformation: To extract the feature of image this algorithm is applied.

4. Results

OCR is used to extract the feature of the image, instead of OCR we are using the SIFT algorithm to extract the features of handwritten document's image. After applying the algorithms for filtering the image, then SIFT Algorithm will be applied, the features of new image will be extracted and matched with saved image's feature.

5. Application areas

This can be use in any application where there is need of Identifying the writer of handwritten document, but specially this will use in

1. In Banking application for signature verification
2. The system it will be use in Criminal case to identify writer.

6. Conclusion

In this paper, we have proposed a SIFT Algorithm which will be used for feature extraction and the feature matching of the Handwritten document. We will estimate the parameters using structured learning and considering the segmentation problem as a binary quadratic programming. By using this function to the proposed formulation, we concentrate on unary properties of gap and pair-wise similarities between word-separator in the word segmentation. By using the Structured Structure Vector Machine, all parameters will estimated and stored in feature vector

References

- Izadi, S., Haji, M., & Suen, C. Y. (2008). A new segmentation algorithm for online handwritten word recognition in Persian script. In *Proc. Eleventh International Conf. Frontiers in Handwriting Recognition (CFHR 2008)* (pp. 598-603).
- Koo, H. I., & Cho, N. I. (2012). Text-line extraction in handwritten chinese documents based on an energy minimization framework. *Image Processing, IEEE Transactions on*, 21(3), 1169-1175.
- Lazimy, R. (1982). Mixed-integer quadratic programming. *Mathematical Programming*, 22(1), 332-349.
- Papavassiliou, V., Stafylakis, T., Katsouros, V., & Carayannis, G. (2010). Handwritten document image segmentation into text lines and words. *Pattern Recognition*, 43(1), 369-377.
- Papavassiliou, V., Stafylakis, T., Katsouros, V., & Carayannis, G. (2010). Handwritten document image segmentation into text lines and words. *Pattern Recognition*, 43(1), 369-377.

Ryu, J., Koo, H. I., & Cho, N. I. (2014). Language-independent text-line extraction algorithm for handwritten documents. *Signal Processing Letters, IEEE*, 21(9), 1115-1119.

Ryu, J., Koo, H. I., & Cho, N. I. (2015). Word segmentation method for handwritten documents based on structured learning. *Signal Processing Letters, IEEE*, 22(8), 1161-1165.

Salvi, D., Zhou, J., Waggoner, J., & Wang, S. (2013, January). Handwritten text segmentation using average longest path algorithm. In *Applications of Computer Vision (WACV), 2013 IEEE Workshop on* (pp. 505-512). IEEE.