

A Novel Approach to Detect and localize the Text in Natural scene image

Prof. B. K. Sarojini, Mr. Mahesh Mahadar

Department of Electronics and Communication Engineering, Basaveshwara Engineering College, Bagalkot India

Abstract: The information available in the natural scene provides very important clues for many image based application. So the detection and localization of text from natural scene image is important task for content based image analysis. In this paper, we develop a novel approach to detect and localize text in natural scene images. This problem is challenging due to the fact that text has different in size, alignment, orientation, style, complex background of images as well as image having the low contrast. Initially, the RGB image is converted into grayscale image and applied to the local binarization algorithm for the segmentation. After that, Haar Wavelet Transform (DWT) is used which is fast as compared with all wavelets because its filter coefficients are either 1 or -1. It decompose image into four sub band, one is average and other three are detailed which helps to find out the approximately confident region. To filter out the non-text component, a conditional random field (CRF) is used which applied on confident region image. Finally, by using some predefined condition, the text is obtained in bounding box.

Keywords: Wavelet transform, Text Detection, Text Localization, Conditional random field (CRF).

I. INTRODUCTION

Recently the digital image capturing devices, such as digital cameras, mobile phones are increased to the large extent so content based image analysis techniques are receiving intensive attention in recent years. As compared with the other contents in images, text information has great advantage because it is easily understood by both human and computer and finds wide range of applications but localization of text is very difficult task for the following reasons. First of all, text has wide range of variety in font, orientation; style as well as size may change from small to too large. Secondly, text present in an image may have low contrast, multiple colors and appear in a much cluttered background. [5]

The text detection and localization methods can be mainly categorized into two groups: connected component (CC) based and region based. Region based methods use texture analysis principle to detect and localize text regions. It extracts the feature vector from each local region and fed into a classifier which estimates the text component. Then neighboring text regions are grouped to generate text blocks. As the text regions have distinct textural properties from non-text ones, these methods can easily detect and localize text even when images are noisy. On the other hand, CC based methods use the edge detection or color clustering for the direct segment and identify text components. Finally, the non-text components are then removed with some specific rules but this method has lower computation cost and the located text components can be used as it is for recognition. Still existing methods have some problems to be solved. Such as for region based methods, they are very slow and the performance is also sensitive to text alignment and orientation. For CC based methods it required a prior knowledge of text position and scale to segment text components accurately. [8]

II. RELATED WORK

Most of the region based methods are mainly based on principle that text regions have different characteristics from non-text regions such as the distribution of gradient strength and texture properties. Many efforts have been made for text extraction and recognition in image.

Weinman et al. provides a method which uses a Conditional Random Field (CRF) model for text detection. This model assigns the candidate components to one of the two classes such as "text" or "non-text" by considering the properties of unary components as well as relationship between the contextual components. This method provides the benefit over the traditional local region based text detection methods. [1]

Chung Wei Liang and Po Yueh Chen provide useful and effective approach to extract the text region from static image. They use Haar Discrete Wavelet Transform (DWT) to localize the text region along with the morphological operator to detect edges of candidate text region which is used for the isolation of text data from documented video image. [2]

The method proposed by Kim et al. uses the support vector machine (SVM) to analyze textural properties of text. The intensity of raw pixel which required for analysis of textual pattern are fed to the SVM then continuously adaptive mean-shift (CAMSHIFT) is applied to the result of texture analysis to identify the text region so the combination of SVM and CAMSHIFT provides a robust and efficient approach for text detection. [3]

X. L. Chen et al. provides the combined approach of multiscale edge detection and multiresolution, color analysis, adaptive searching as well as affine rectification in the proposed framework for sign detection, with different emphases at each phase to handle the text in different orientations, sizes, color distributions and backgrounds. It uses affine rectification to recover deformation of the text regions caused by an inappropriate camera view angle. The procedure can significantly improve text detection rate and optical character recognition (OCR) accuracy. [5]

Kim segments an image using color clustering in a color histogram of RGB space. Non-text components mainly contain the long horizontal lines as well as image boundaries these are eliminated by iterative projection profile analysis. To filter out the non-character components this method used cluster-based templates for multi-segment characters to lower down the difficulty in defining heuristics for filtering out non-text components. [4]

Zhanget al. uses the Markov random field (MRF) to detect the neighboring information of components. First mean shift algorithm is used to segment the candidate text components then after developing the component adjacency graph, first-order component term are integrated using MRF model and finally for labeling components as "text" or "non-text" a higher order contextual term is used. [7]

S. Audithan et al. uses the Haar discrete wavelet transform (DWT) which is the fastest among all wavelets because its filter coefficients are either 1 or -1. DWT detected edges and then based on the edge map, line feature vector graph is generated which helps to extract the stroke information. Finally text regions are generated and filtered according to obtain the line features. [10]

Zhu et al. uses the first nonlinear local binarization algorithm for the segmentation of candidate CCs. There are large numbers of component feature which contains geometry, edge contrast, shape regularity, spatial coherence features as well as stroke statistics are defined to train an AdaBoost classifier which helps to filter out non-text components. [6]

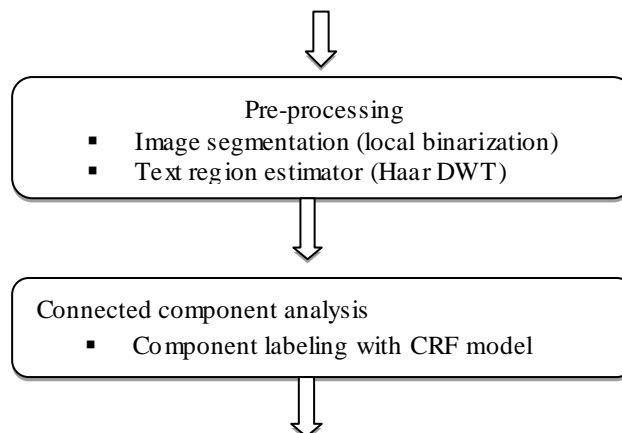
The approach of Khushbu C. Saner is mainly divided into three steps 1) Pre-processing 2) Connected Component Analysis 3) Optical Character Recognition. In pre-processing step, the color image is converted into binarized image for the edge detection. A conditional random field (CRF) model use binary discourse part relationships and unary part properties along with the designed supervised parameter learning to filter out the non-text components. Then recognized text is localized in original image and text parts are classified into text lines. Once the line partition is over, character recognition will be done using Optical character recognition to recognize the character. Here the results are evaluated on the natural image dataset. [9]

Liu et al. provide the approach to detect color texts from natural scene images. It contains the combination of connected component based approach as well as region based approach. To detect the probabilities of text scale and position a text region detector is designed then an efficient local binarization algorithm is used to segment candidate text components. A conditional random field (CRF) model along with supervised parameter learning is designed to combine the binary contextual component relationships and unary component properties. Finally, learning-based energy minimization method is used to group the text components into text lines or words. [8]

III. SYSTEM OVERVIEW

To overcome the above difficulties, we present an approach to detect and localize texts in natural scene images with wavelet transform. A text region detector is designed using Haar discrete wavelet transform to estimate the probabilities of text position, and then segment candidate text components using an efficient local binarization algorithm. [6] For labeling the connected components as a text or non-text, we use the unary as well as binary properties of the conditional random field. Finally, text components are localizing with use of bounding box. [1] Figure 1 shows the flowchart of proposed system.

Input image



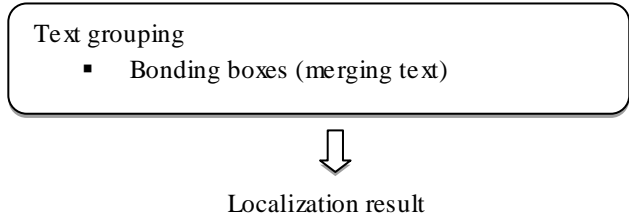


Figure1: Flow Chart of Proposed Model

VI. PRE PROCESSING

A text region detector is designed to estimate the text confidence and the corresponding scale which helps to efficiently utilize and extract local text region information, depends on which candidate text components can be segmented and analyzed accurately.

A. Image Segmentation

Initially RGB image is converted into gray-level image. Then Niblack's local binarization algorithm is used because of its high insensitivity and efficiency for degraded image. To segment candidate connected components (CCs) from the gray-level image, the following formula should be used to binarize each pixel which is defined as

$$m(x) = \begin{cases} 0, & \text{if } \text{gray}(x) < i(x) - k \cdot s(x); \\ 255, & \text{if } \text{gray}(x) > i(x) + k \cdot s(x); \\ 100, & \text{otherwise,} \end{cases}$$

Where $s(x)$ and $i(x)$ are the standard deviation (STD) and intensity mean of the pixels for a radius window which has pixel x is available at the center and k is smoothing term which empirically set to 0.45. Most of the methods use the window having fixed radius or it is chosen based on some simple rules such as the gray-level standard deviation (STD) while in our method we use text scale map to calculate the radius of window which provides more stability under noisy conditions. In local binarization, we assume that within each local region, the foreground pixels has gray-level values must be higher or lower than the average intensity of the pixels, So the connected components with a value 0 or 255 are extracted as candidate text components while those of value 100 are non-text components so they are not considered further. [6] Figure 2 shows the resultant image of image segmentation.

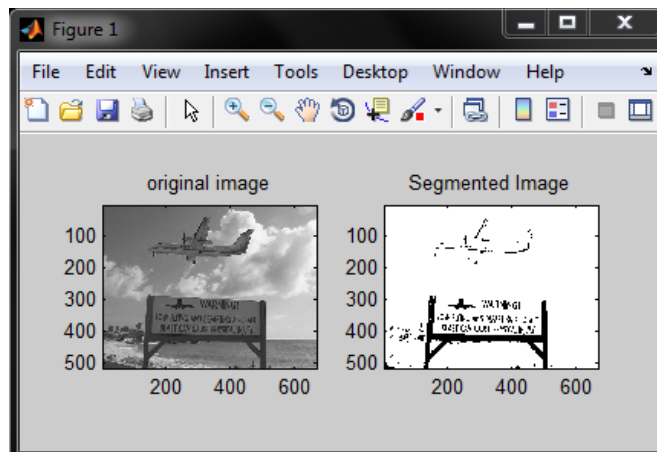


Figure 2: (a) original image (b) Segmented image

A. Text Region Estimator

If the input image is a gray-level image, such image is directly processed to Haar discrete wavelet transform for text region estimation. But if the input image is colored, then its intensity image should be calculated by combining its RGB components. Normally, color images are captured with the help of digital cameras. These pictures are mainly available in the Red Green Blue color space. Intensity image Y is calculated as:

$$Y = 0.299R + 0.587G + 0.114B$$

These image Y is then finally processed with 2-D Haar discrete wavelet transform. Here the Y actually represent the Value component of the Hue Saturation Value (HSV) color space. So there is conversion from RGB color space into HSV color space in above step, once it is over then the Value component is extracted from HSV color space using above expression. To reduce the effect of noise in the image mostly median filtering techniques are used which is applied on the above grayscale image. After this filtering step, a major part of noise will be removed while the edges in the image are still preserved. [2]

B. Haar discrete wavelet transforms

For multi resolution representation, a Haar discrete wavelet transform is a very powerful tool for image processing as well as signal analysis. It can decompose signal in the frequency domain with different frequency components. Two dimensional discrete wavelet transform decomposes an input image into four components, one average component (LL) and three detail components (LH, HL, HH) by calculations of low-pass and a high-pass filter combination as shown in Figure 3.

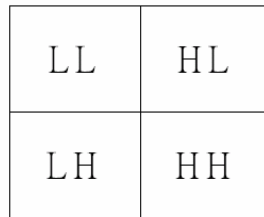


Figure 3: The result of 2-D DWT decomposition

To detect the candidate text edges in the original image the detailed sub-band component are used. In image processing, the multi resolution of 2-D DWT is mainly used to detect edges from the original image. The processing time of the 2-D DWT is much faster than traditional edge detection filters because it can detect three kinds of edges at a time. But the result provided by 2-D DWT can be similar as compared with the conventional edge detection filters. This is why we choose Haar DWT because it is efficient and simpler than that of any other wavelets. [10] Figure 4 shows the approximate text area in image.

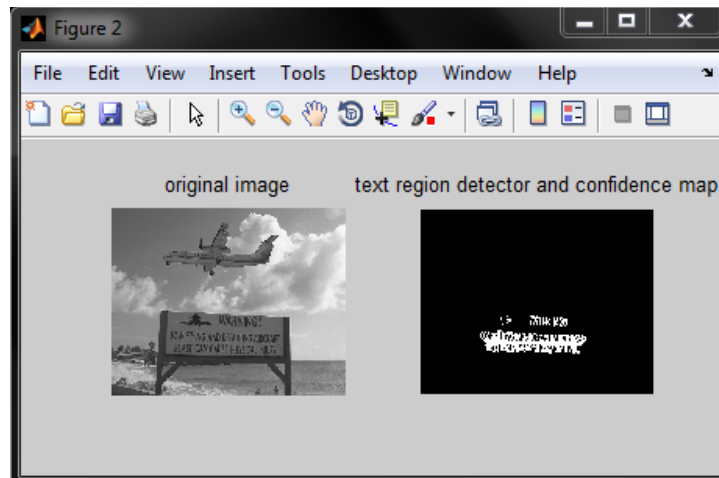


Figure 4: (a) original image (b) text region detector and Confidence map

V. CONNECTED COMPONENT ANALYSIS

For connected components analysis (CCA), we assign each candidate components to one of the two classes such as “text” and “non-text” by using both binary contextual component relationships and unary component properties of conditional random field (CRF) model.

A. Introduction of CRF

Conditional random fields (CRFs) are mainly designed for labeling tasks such as document image segmentation, text identification as well as natural language processing. CRFs are basically probabilistic graphical models and used to label the text region from different areas in images having the spatial interdependencies. For example, text blocks are sequentially available from left to right. So by considering the information of neighboring text blocks, isolated noises available in these sequential text blocks can be easily removed which provide more accurate results of labeling. [1] So CRFs provide a more flexible formulation rather than the other generative graphical models such as Markov random fields (MRFs) which require specifying the likelihood function. [7] More formally, let assume the observed features from candidate blocks $X = \{a_i\}$, and random variables over corresponding labels $Y = \{b_i\}$. Then joint distribution over the label b_i with given observation a_i is represented as

$$d(b_i | a_i) \propto \exp \left(\lambda P(b_i, X) + \mu \sum_{(i,j) \in E} R(b_i, b_j, X) \right) \quad (1)$$

Where the function $P(b_i, X)$ is called as associated potential which measures the confidence of label b_i by considering the observations, and function $R(b_i, b_j, X)$ is interaction potential which provides smooth labels over entire graph G , and λ and μ parameters helps to control the influence from neighboring nodes to center node i and observations, and $(i,j) \in E$ represents the neighboring nodes of node R which are connected by edges E in the graph G .

In our work, we use the topology for our CRFs. By considering a Markov assumption, detected block in the image are exclusively represent by each gray node g_i in the hidden layer and then connect to its four nearest neighbors block along with their corresponding observations. In the case of real images, the neighbor blocks are mainly determined by Euclidean distance between them but it may not necessarily be located as a grid. To integrate the predicted confidence of blocks into CRFs framework, we define the associated potential as

$$P(b_i, X) = \sum_{j \in N} e_j \exp(-|s_{i,j} \cdot \cos(\theta_{i,j})|) \quad (2)$$

Where j runs over neighbors of node i including itself, and S_{ij} is the spatial Euclidean distance between node i, j and e_j is the posterior which is estimated by the SVM for node i, j and θ is the angle between centers of node i and j . The idea behind equation 2 is that if two neighboring nodes are very close to each other as well as their separation is mostly horizontal, then they have more influence on each other. [9]

B. Properties of CRF

a. Unary Component Features

Here we consider different types of unary component features such as Aspect ratio, height, normalized width and Compactness, To characterize single component's geometric and textural properties. [8]

b. Binary Component Features

Here we consider different types of binary component features such as Overlap ratio, Shape difference, Gray-level difference, Scale ratio, To characterize the geometric and spatial relationship and textural similarity between two neighboring

component. [8]

VI. TEXT GROUPING

Here adjacent letters are grouped together to form words. For the performance analysis of a text extraction algorithm, it is recommended that the recall rates and precision must be computed. But the performance parameters are mainly dependent on correctly classified words. Some proposed methods are effective but it is too complicated because of training data necessity while the other methods are simpler but not effective for text grouping. To overcome all such drawbacks we propose the bounding boxes (BB) concept to merge adjacent letters in words which is based on the computation of distances between these boxes (BB) of letters, detected in the above step. The parameters B1 and B2 are used in the merging letters process which represents the center coordinates of the two BBs of connected component. Figure 4 represents the merging process.

Here B1 (y1a), B1 (y2a) and B2 (y1a), B2 (y2a) represent the coordinates of the first and second BBs in vertical direction respectively while Width 1 and Width 2 represent the width of the first and second BBs respectively. 'Distance' represents the distance between the centroid of the two BBs considered in the horizontal direction.

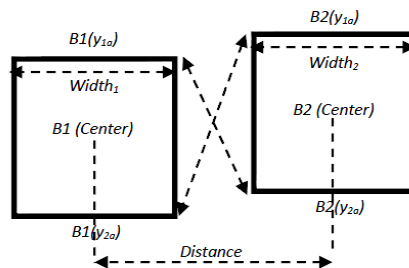


Figure 4: Parameters used in merging process

The first step for merging is based on merging of letters along the horizontal line. Here we consider only those images which contain relatively well aligned letters. The conditions for merging the letters in the detected regions are defined as follows

- $[B2(y2a) > B1(y1a)] \& [B2(y1a) < B1(y2a)]$
- $[Distance < 0.7 \times \text{Max}(Width1, Width2)]$

The pair of BBs which satisfy both the above conditions is then merged together in this step to obtain the word. [8] Figure 5 shows the bounding box on each text in input color image.

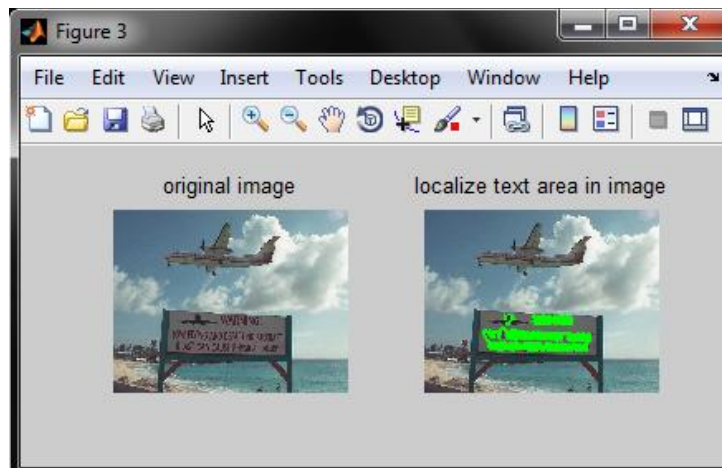


Figure 5: (a) Original image (b) localize text area in image

VII. CONCLUSION

The given input color image to be converted into grayscale image and then image segmentation is carried out on the grayscale image to obtain the segmented image. Wavelet Transform used to decompose the segmented image. It will decompose the original image into four frequency sub bands to improve the contrast and resolution of the image which helps to find the approximated text area in the image using connected component graph. Finally by using bounding box concept, we localize text in image. So it provides robust approach to detect and localize texts by integrating region information.

REFERENCES

1. J. Weinman, A. Hanson, and A. McCallum, "Signal detection in natural images with conditional random fields," in Proc. 14th IEEE Workshop on Machine Learning for Signal Processing (MLSP'04), São Luis, Brazil, 2004, pp. 549-558.
2. Chung-Wei Liang and Po-Yueh Chen, "DWT based Text Localization", International Journal of Applied Science and Engineering: 2004
3. K.I. Kim, K. Jung, and J.H. Kim, "Texture-based approach for text detection in images using support vector machines and continuously adaptive mean shift algorithm," IEEE Trans. Pattern Anal. Mach. Intell. vol. 25, no. 12, pp. 1631-1639, 2003.
4. K. Jung, K.I. Kim, and A.K. Jain, "Text information extraction in images and video: A survey," Pattern Recogn., vol. 37, no. 5, pp. 977-997, 2004.
5. X. L. Chen, J. Yang, J. Zhang, and A. Waibel, "Automatic detection and recognition of signs from natural scenes," Jan. 2004.
6. K.H. Zhu, F.H. Qi, R.J. Jiang, L. Xu, M. Kimachi, Y. Wu, and T. Aizawa, "Using AdaBoost to detect and segment characters from natural scenes," in Proc. 1st Conf. Camera Based Document Analysis and Recognition (CBDAR'05), Seoul, South Korea, 2005, pp. 52-59.
7. Zhang and S.F. Chang, "Learning to detect scene text using a higher order MRF with belief propagation," in Proc. IEEE Conf. Computer Vision and Pattern Recognition Workshops (CVPRW'04), Washington, DC, 2004, pp. 101-108.
8. Yi-Feng Pan, Xinwen Hou, and Cheng-Lin Liu, Senior Member, IEEE, "A Hybrid Approach to Detect and Localize Texts in Natural Scene Images" march 2011
9. Khushbu C. Saner, "Robust Approach to Recognize and Localize Text from Natural Scene Images," Oct 2014.
10. S. Audithan, R.M. Chandrasekaran, "Document Text Extraction from Document Images Using Haar Discrete Wavelet Transform," Vol. 36 No. 4, pp. 502-512, 2009.