

# A Robust Algorithm for Text Detection in Color Images

Yangxing LIU, Satoshi GOTO, Takeshi IKENAGA

## Abstract

*Text detection in color images has become an active research area since recent decades. In this paper, we present a novel approach to accurately detect text in color images possibly with a complex background. First, we use an elaborate edge detection algorithm to extract all possible text edge pixels. Second connected component analysis is employed to construct text candidate region and classify part non-text regions. Third each text candidate region is verified with texture features derived from wavelet domain. Finally, the Expectation maximization algorithm is introduced to binarize text regions to prepare data for recognition. In contrast to previous approach, our algorithm combines both the efficiency of connected component based method and robustness of texture based analysis. Experimental results show that our algorithm is robust in text detection with respect to different character size, orientation, color and language and can provide reliable text binarization result.*

## 1. Introduction

The retrieval of text information from color images has gained increasing attention in recent years. Text appearing in images can provide very useful semantic information and may be a good key to describe the image content. Many papers about the text detection from color images or video sequence have been published. However, due to the complexity of text appearance in color images, text detection is still a difficult and challenging task in image processing.

Problems with many existing methods are that they can not perform well in the case of variant text orientation, size and low resolution image, where characters may be touching. Many methods assume that the text direction is horizontal or vertical [1] and text font size is in limited range. Some proposed methods [2] fail to detect isolated character because there is no contextual information and certain text parameters, such as base line, character width and height, can not be accumulated with statistical method in such case.

In this paper, a novel text detection method targeted towards being robust with respect to diverse kinds of

text appearances, including character font size, orientation, color and language, is presented. First, with an elaborate edge detection algorithm, we extract the edge pixels of all possible text regions in color images. Then connected components generated from the edge map are carefully examined to yield text region candidates. This is followed by the analysis of texture features of candidate text regions, instead of the whole image, to verify true text regions and separate non-text regions. Eventually, the Expectation maximization (EM) algorithm is introduced to binarize each text region to prepare data for OCR (Optical Character Recognition). Since our algorithm is applied to the edge maps, which are more stable than complete color images under varying illumination conditions, and not to entire image, it is robust to intensity variations.

The rest of this paper is organized as follows. Section 2 gives an overview of previous work on text detection. Section 3 describes the proposed algorithm in detail. Experimental results are explained in section 4. Conclusion remarks are given in last section.

## 2. Related work

Many methods for text detection in images and videos have been published in the past. Based on the ways being used to locate text regions, most of the text detection techniques can be classified as either connected component(CC) or texture based algorithms.

The first class is based on the analysis of the geometrical arrangement of edges or homogeneous color and grayscale components that belong to characters. The CC based algorithms are relatively simple to implement, but they are not very robust for text localization in images with complex background. Chen et al [3] detected vertical and horizontal edges in an image and dilated two kinds of edges using different dilation operators. Real text regions are then identified by using support vector machine. Zhong et al [4] extracted text as those connected components that follow certain size constraints and horizontal alignment constraints. Jain and Yu [5] also use connected component analysis to localize text candidates. But this method can extract only horizontal texts of large sizes.

Texture-based methods are based on the fact that texts in images have distinct textural properties that can be used to discriminate them from the background or other non-text region. In general, the textured-based algorithms are more robust than the connected component-based algorithms in dealing with complex background [6]. The high complexity of texture segmentation is the main drawback of this method. The algorithm proposed in [1] using texture features to extract text but failed in the case of small font characters. In [7], wavelet features from fix-size blocks of pixels and classify the feature vectors into text or non-text using neural networks. Because the neural network based classification is performed in the whole image, the detection system is not very efficient in terms of computation cost.

### 3. Text detection algorithm

In contrast to previous approaches, our algorithm is a hybrid approach, which combines connected component based and texture based method. First with image edge detection result, we generate connected component of text region candidate efficiently. Then texture analysis in wavelet domain is performed on each text region candidate to verify the text region and remove false alarms (i.e. non-text regions).

In the following, we explain and state the details of four main steps of our algorithm in order of processing.

#### 3.1. Edge Detection

In this stage, we should extract precise edge information of all text regions and filter out edges of most non-text regions in image. An important fact we observed is that the pixels representing text contour usually have a high contrast to their neighbor pixels. In our algorithm we use black pixels to represent the edge pixels and white to represent non-edge pixels.

Because the edge detector we adopted is different from most existing isotropic edge detector, which can not provide accurate edge direction information, basic ideas are first presented below to illustrate the edge detector. Some ideas about the detector have been described in [8].

##### 3.1.1. Edge Detector

Given a color image, the difference vector DV in rgb color space induced by moving an infinitesimal step in the image plane in the direction  $\{dx, dy\}$  is :

$$DV = (dx \ dy) J_c^T,$$

$$J_c = \begin{bmatrix} \partial r / \partial x & \partial r / \partial y \\ \partial g / \partial x & \partial g / \partial y \\ \partial b / \partial x & \partial b / \partial y \end{bmatrix} = \begin{bmatrix} r_x & r_y \\ g_x & g_y \\ b_x & b_y \end{bmatrix}$$

where  $J_c$  is the Jacobian matrix of the image.

The Euclidean squared magnitude of DV is  $DV^2 = (dx \ dy) M_c (dx \ dy)^T$ ,

$$\text{where } M_c = J_c^T J_c = \begin{bmatrix} M_{xx} & M_{xy} \\ M_{xy} & M_{yy} \end{bmatrix},$$

$$M_{xx} = (r_x)^2 + (g_x)^2 + (b_x)^2,$$

$$M_{xy} = r_x \times r_y + g_x \times g_y + b_x \times b_y,$$

$$M_{yy} = (r_y)^2 + (g_y)^2 + (b_y)^2,$$

and asking for the direction of  $\{dx, dy\}$  maximizing this magnitude is an eigenvalue problem. We can obtain the magnitude extremum in the direction of the eigenvector of the matrix  $M_c$  and the extremum value is the corresponding eigenvalue. So we can get precise gradient magnitude and direction of each image pixel by computing corresponding eigenvalue and eigenvector direction.

##### 3.1.2. Edge Extraction Algorithm

First, median filter was applied to the input image C to reduce the noise of the input images while preserving sharp edges.

Then the matrix  $M_c$  is computed for each pixel on the image and we will get a series of eigenvalues V and eigenvectors E that reflect the gradient magnitude and direction of each pixel.

As we know, edge pixels are those points with local maximum gradient magnitude in their gradient direction. Furthermore, we also notice that the edges of text symbols are typically stronger than those of noise or background areas. So a pixel (i,j) is accepted as an edge pixel only if it meets the following two requirements.

- First, the pixel must have a larger gradient magnitude than that of its neighbor located in the direction closest to its gradient direction, i.e., the eigenvalue of an edge pixel must be greater than that of both two neighboring pixels, which are closest to its eigenvector direction.
- Second, the pixel gradient magnitude (i.e. eigenvalue  $V(i,j)$ ) must be greater than the adaptive threshold T to eliminate weak edges. T is determined by the following formula:

$$T = \frac{\sum_{(i,j) \in C} (V(i,j) \times |V_{dif}(i,j)|)}{\sum_{(i,j) \in C} |V_{dif}(i,j)|},$$

where  $V_{\text{dif}}$  is the corresponding eigenvalue difference between two neighboring pixels closest to the eigenvector direction of pixel  $(i,j)$ .

To make character edge more continuous, conditional dilation is performed on edge collection obtained in previous operation to connect text edges to form closed contours. A  $3 \times 3$  square structuring element with the origin at its center is selected to dilate the image edge. Furthermore, the gradient magnitude of center pixel must exceed  $T$  and the gradient direction (corresponding eigenvector direction) difference between the center pixel and its neighboring edge pixel must be less than a constant  $\alpha$ , which has been experimentally set to 0.26, about 15 degree.

After this, all character edge pixels as well as some non-character edge pixels which also show high local color contrast are remained in the image edge map. Then we can link connected edge pixels through connected component analysis to generate candidate text regions.

### 3.2. Candidate text regions generation

In this phrase, we first need to group connected edge pixels into different regions.

Unlike previous work, we extract connected component from text contour. Our connected component analysis is performed on black edge pixels generated in the previous step. A connected component is defined as a set of black pixels where each pixel is a direct neighbor of at least one other black pixel in the component. Our connected component generation basically follows the 8-neighborhood-connectivity algorithm discussed below.

Scan the image from left to right and from top to bottom. Initialize the class label of each edge pixel  $CL(i,j)$  to number 0. Given an edge pixel  $P$ , examine the four neighbors of  $P$ , which have already been encountered in the scan (i.e. the neighbors to the left of  $P$ , above it, and the two upper diagonal terms). Then we can label  $P$  according to following different cases.

- If the class number of all four neighbors are 0, assign a new class number to  $P$ , else
- If only one neighbor has class number bigger than 0 and the gradient direction difference between this neighbor and  $P$  is less than  $\alpha$ , assign this class number to  $P$ , else
- If more than one neighbor has class number bigger than 0 and the gradient direction difference among these neighbor pixels and  $P$  is less than  $\alpha$ , assign one of the labels to  $P$  and make a note of the equivalences, else
- Assign a new class number to  $P$ .

After completing the scan of the whole image, the equivalent label pairs are sorted into equivalence classes and a unique label is assigned to each class. Then a second scan is made through the whole image, during which each label is replaced by the label assigned to its equivalence classes.

Using above algorithm, edges of a character can be linked into a closed contour and we may obtain many connected components from edge map. The next stage is selecting part of the components as text region candidate and removing false alarms.

To develop the criteria for filtering out non-text regions, the following features of text are observed:

#### 1) Average edge pixel gradient magnitude

The first fact is that the average edge pixel gradient magnitude is usually higher for text than non-text blocks. We can calculate this average value of each connected component region  $R$  with following formula:

$$V_{\text{avg}} = \frac{\sum_{(i,j) \in R} V(i,j)}{m}, (CL(i,j) = CL_R),$$

where  $m$  is the number of pixels labeled with one connected component class  $CL_R$ . The  $V_{\text{avg}}$  of a text region should be greater than  $2 \cdot T$ .

#### 2) Edge gradient direction variance

Another important fact we noticed is that text regions have a higher gradient direction distribution variance than graphic regions. The variance of a region can be estimated by  $\theta_{\text{max}} - \theta_{\text{min}}$ , where  $\theta_{\text{max}}$  and  $\theta_{\text{min}}$  are the maximum and minimum edge gradient direction angle respectively in a connected component region. Within a text region, the variance must be greater than  $\pi$  (180 degree), i.e.  $(\theta_{\text{max}} - \theta_{\text{min}}) > \pi$ .

#### 3) Edge pixel number

The third valuable fact we observed that text block should have more edge pixels than some non-text block. The edge pixel count in a text block, labeled as the same class within a connected component region, must be greater than  $\text{MAX}(2 \cdot W, 2 \cdot H)$ , where  $W$  and  $H$  are width and height of the connected component region respectively.

With features described above, three criteria are applied on every connected component in order to reject blocks constructed of noise pixels and classify some non-text blocks.

### 3.3. Text region verification

The text region localization described in previous section may produce false alarm, which also have high edge density and strength. So we need to perform more accurate texture analysis to classify text regions.

Instead of performing a global texture analysis on the whole image, we consider every region of interest separately to remove false alarms, which makes our algorithm more efficient and robust.

Text shows a rhythmic spatial pattern distribution, which consists of a regular alternation of contrast changes in one specific direction. Therefore texture analysis can be exploited to separate non text regions like graphics, images and other non-text rectangular homogeneous and high contrast regions.

We select the Harr wavelet as the basis for texture characterization because of its good ability to characterize texture features and its computation efficiency [9].

Let  $\phi(x)$  and  $\psi(x)$  be the two Haar wavelet bases of one-dimensional.

$$\begin{cases} \phi_{k,s}(x) = 2^{-k/2} \phi(2^{-k}x - s) \\ \psi_{k,t}(x) = 2^{-k/2} \psi(2^{-k}x - t) \end{cases}$$

With 2-dimensional image data, the corresponding tensor product transform basis can be calculated as follows:

$$\begin{cases} \phi_{k,s,t}^{LL}(i,j) = \phi_{k,s}(i)\phi_{k,t}(j), \phi_{k,s,t}^{HL}(i,j) = \phi_{k,s}(i)\psi_{k,t}(j) \\ \phi_{k,s,t}^{LH}(i,j) = \psi_{k,s}(i)\phi_{k,t}(j), \phi_{k,s,t}^{HH}(i,j) = \psi_{k,s}(i)\psi_{k,t}(j) \end{cases}$$

Because the input image is colored, first the input image is converted into grey-level image I. Then the image I is processed with discrete wavelet transform and transformed into four sub-bands LL, HL, LH and HH with a 2-channel filter bank (L: low pass filter, H: high pass filter).

After decomposing the image I into 2-D Haar wavelet, we can compute wavelet moment features to capture each candidate text region texture property.

The wavelet energy feature of a pixel (i,j) is defined as:

$$\text{ENG}(i,j) = |\text{LH}(i,j)| + |\text{HL}(i,j)| + |\text{HH}(i,j)|.$$

Given a text region candidate R with  $N_e$  edge pixels labeled with  $\text{CL}_R$ , the third order wavelet moment  $M_3(R)$  can be calculated as

$$M_3(R) = \frac{1}{N_e} \sum_{(i,j) \in R} |(\text{ENG}(i,j) - \text{MENG}(R))^3| \cdot (\text{CL}(i,j) = \text{CL}_R),$$

where  $\text{MENG}(R)$  is the mean energy feature of region R and

$$\text{MENG}(R) = \frac{1}{N_e} \sum_{(i,j) \in R} \text{ENG}(i,j), (\text{CL}(i,j) = \text{CL}_R).$$

For each text region candidate, its third order wavelet moment is checked to verify whether it is a text region or not. The candidate region R is rejected as non-text region if its third moment value is greater than  $T_m$ , which can be calculated with equation (1).

$$T_m = T^2 * N_e * 2.5 \quad (1)$$

### 3.4. Text region binarization

After we verify each candidate text region with its texture property, we can binarize each text region separately and feed the result to OCR.

Although human may perceive single character with the same color appearance, the actual pixel colors may vary significantly. So it is necessary to perform color clustering to compensate for these effects.

In each text region, we just want to cluster all colors into two distinctive colors to discriminate between text and other non-text part. So we use a mixture model of Gaussians described with equation (2) to depict the color distribution in text regions, where  $x$  is an (r, g, b) vector. Because EM algorithm is best suited for fitting Gaussian clusters, we use it to estimate the distribution parameter. The EM algorithm iterates by adjusting the parameters of the Gaussian probability model to maximize the likelihood of data in dataset. The iteration stops when the difference between two successive iterations becomes negligible.

$$P(x) = \sum_{i=1}^n \frac{p(i) \times \exp\left\{-\frac{1}{2}(x - \mu_i)^T \sum_i^{-1}(x - \mu_i)\right\}}{(2\pi)^{(d/2)} \left|\sum_i\right|^{1/2}} \quad (2)$$

where  $\mu_i$  is the mean of cluster  $i$ ,  $\sum_i$  is the covariance matrix of cluster,  $d$  is the dimensionality of the data and  $d = 3$  with color image (rgb component).

In our work, we use the filling algorithm to select two different colors of two pixels inside and outside in text contour to carefully initialize for EM algorithm, (but we have found that the initialization has little effect on the quality of the resulting clustering process from experiment results).

Upon convergence of the EM algorithm, the two mean vectors can be recorded as the two dominant colors in a text region, i.e. text and non-text part color. Thus we can binarize this text region via measuring the Euclidian distance between each pixel color and two mean color vectors.

## 4. Experimental results

Currently, our algorithm has been implemented in C++ language under Windows-XP on an EPSON Endeavor MT7000 PC.

To evaluate the actual performance of our proposed algorithm, we tested 415 real color images which include different types of texts. The image resolution range is between 73\*42 and 3072\*2048 pixels. Among

them, 200 images are from the conference ICDAR 2003 scene text detection competition dataset (529 scene images). Another 215 real color images are carefully chosen with a wide variety of background complexity and text types. Text appearance varies with different colors, orientation and languages and the character font size in images ranges from 8pt to 530pt. The classification result is given in table I.

For measuring accuracy, recall and false alarm rate are calculated to evaluate our algorithm performance. Recall rate is defined as follows,

$$[\text{Recall rate}] = \frac{\text{number of detected characters}}{\text{number of characters in image}} \times 100\% .$$

False alarm rate is evaluated as follows,

$$[\text{False alarm rate}] = \frac{\text{total pixel count in detected non - text regions}}{\text{number of pixels in image}} \times 100\% .$$

Experimental results show that our algorithm has a high recall rate with low false alarm rate. Table I shows evidence that our texture analysis led to small deficits for recall rate but greatly reduced the false alarm rate from 11.5% to 3.7%. Some false alarms left because of their strong texture features, high contrast and resemble text-like attributes. So we plan to combine recognition engine to detection process to further reduce the false alarm in future.

Experimental results also show that the binarization procedure with EM algorithm is also quite usable. Even when the color distribution inside a region is unimodal, such as characters like “l”, “—” etc., the two mean vectors become nearly coincident, which makes the binarization process more robust.

## 5. Conclusions

In this paper, we propose a hybrid approach, which combines connected component based and texture based methods, to detect variant texts in color images. First, an elaborate image edge extraction algorithm and connected component analysis are used to extract candidate text regions. Then texture feature in wavelet domain is explored to identify text regions from candidate regions. After text region verification, EM

algorithm is introduced to binarize all text regions to prepare data for OCR. Experimental results demonstrate that our algorithm is independent from different text orientation, size and language. Furthermore, our algorithm can also detect an isolated character, which has no neighboring text region.

## References

- [1] Kwang In Kim, Keechul Jung and Jin Hyung, “Texture-based approach for text detection in images using support vector machines and continuously adaptive mean shift algorithm”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, issue 12, pp. 1631-1639, Dec. 2003.
- [2] Yefeng Zheng, Huiping Li and Doermann, D., “Text identification in noisy document images using Markov random model”, In *Proceedings of the Seventh International Conference on Document analysis and Recognition*, vol.1, pp. 599-603, Aug. 2003.
- [3] Chen Datong, Bourlard H. and Thiran J. P., “Text identification in complex background using SVM”, In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, pp. 621-626, Dec. 2001.
- [4] Y. Zhong, K. Karu and A.K. Jain. “Locating Text in Complex Color Images”, *Pattern Recognition*, vol. 28, no. 10, pp. 1523-1536, October 1995.
- [5] Jain A.k. and Bin Yu, “Automatic text location in images and video frames”, In *Proceedings of the Fourteenth International Conference on Pattern Recognition*, vol. 2, 16-20, pp. 1497-1499, Aug. 1998.
- [6] Xiaou Tang, Xinbo Gao, Jianzhuang Liu and Hongjiang Zhang, “A spatial-temporal approach for video caption detection and recognition”, *IEEE Transactions on Neural Networks*, vol. 13, issue 4, pp. 961-971, July 2002.
- [7] Fujii, Masafumi., Wolfgang, J. R. and Hoefer, “Filed-Singularity correction in 2-D time-domain Haar-wavelet modeling of waveguide components”, *IEEE Transactions on Microwave Theory and Techniques*, vol. 49, issue 4, pp. 685-691, Apr. 2001.
- [8] Lee, H.-C., Cok, D.R., “Detecting Boundaries in a Vector Field,” *IEEE Transactions on Signal Processing*, vol. 39, issue 5, pp. 1181-1194, May 1991.
- [9] Mojsilovic A., Popovic M.V. and Rackov D.M., “On the selection of an optimal wavelet basis for texture characterization”, *IEEE Transactions on Image Processing*, vol. 9, issue 12, pp. 2043-2050, Dec. 2000

Table I Text Detection Result

	Character Language			Character font size (pixels)			Character Orientation		
	English	Japanese	Chinese	8~200	201~300	301~530	Vertical	Horizontal	Arbitrary
Number of characters	8653	932	542	8270	1131	726	3735	5428	964
	Before Texture Analysis								
Recall rate	91.3%	90.5%	92.7%	91.4%	91.1%	90.5%	90.8%	91.2%	93.8%
False alarm rate	11.5%								
	After Texture Analysis								
Recall rate	91.1%	89.9%	92.1%	91.3%	90.8%	88.4%	90.6%	90.9%	93.5%
False alarm rate	3.7%								