# Document Layout Analysis Based on Fuzzy Energy Matrix

**KangHan Oh, SooHyung Kim***

School of Electronics and Computer Engineering
Chonnam National University, 77 Yongbong-ro, Buk-gu, Gwangju 500-757, Korea

## ABSTRACT

*In this paper, we describe a novel method for document layout analysis that is based on a Fuzzy Energy Matrix (FEM). A FEM is a two-dimensional matrix that contains the likelihood of text and non-text and is generated through the use of Fuzzy theory. The key idea is to define an Energy map for the document to categorize text and non-text. The proposed mechanism is designed for execution with a low-resolution document image, and hence our method has a fast processing speed. The proposed method has been tested on public ICDAR 2009 datasets to conduct a comparison against other state-of-the-art methods, and it was also tested with Korean documents. The results of the experiment indicate that this scheme achieves superior segmentation accuracy, in terms of both precision and recall, and also requires less time for computation than other state-of-the-art document image analysis methods.*

*Key words: Fuzzy Energy Matrix, Document Layout Segmentation, Fuzzy Set.*

## 1. INTRODUCTION

In recently years, high resolution phone cameras have been developed and are widely used in the smart phones to get better image. The document recognition based on smart phone is also attractive theme for smart phone user. Therefore new application focus on algorithm time consuming in order to lunch smart phone environment. In the document recognition field, document image layout analysis is necessary preprocessing part. Usually, a scenario of document recognition consists of three tasks. First, a document layout analysis application segments text and non-text. Second, text lines and words are categorized. Finally, the OCR engine recognizes word information. Therefore layout analysis algorithm should be fast and precise for next part. However traditional document analysis approaches [6]-[10] tend to undervalue time consuming and made algorithm heuristic for increasing performance. So our algorithm considers not only performance but also processing time.

In the past few decades various document layout analysis method have been proposed using different paradigms. There are generally two models for documents analysis: bottom-up and top-down. The bottom-up approach relies on low-level features such as pixel and contented component result. In the document analysis, generally, categorized components can be merged into successively group of blocks represent paragraph or figure [6]-[9], [11]. Whereas, the top-down approaches divide page into group of blocks, which are successively split

into sub-blocks, in an iterative mechanism, to obtain the text, figure [10]. Usually, top-down method is faster than bottom-up but it is less accurate than bottom-up.

In the bottom-up approaches, Koich at al [9] proposed a document analysis algorithm using Voronoi diagram. In order to reduce processing time, they focus on text segmentation because quality of segmented non-text is not important in the OCR engine process for recognition. Therefore our method also does not spend much time for accurately detecting non-text. Laura et al [8] utilized a neuro-fuzzy methodology to segment document. The proposed method strategy is capable of describing the physical structure of document. They extract pixel level features using fuzzy set and then after training, an adaptive neuro fuzzy network algorithm recognizes text and non-text. In [11] used top-down approach based on the connected component extraction. The proposed document analysis model, which preserves top-down generation, is proposed based on which a document is logically represented for interactive editing, storage, retrieval, transfer, and logical analysis. The authors have demonstrated pretty good results. Henry at el [10] extracts background region in the document image using computational-geometry algorithms for off-line enumeration of maximal white rectangles and on-line rectangle unification. This algorithm is simple and fast however algorithm is so heuristic. In the limited datasets, heuristic algorithm can reach to high performance but it cannot expect good performance in different environment.

In our work, the fuzzy theory with rectangular membership function is used for just converting continuous input values to several grade values in order to make FEM. The reason for taking fuzzy theory is that the continuous input values are not suitable for generating FEM.
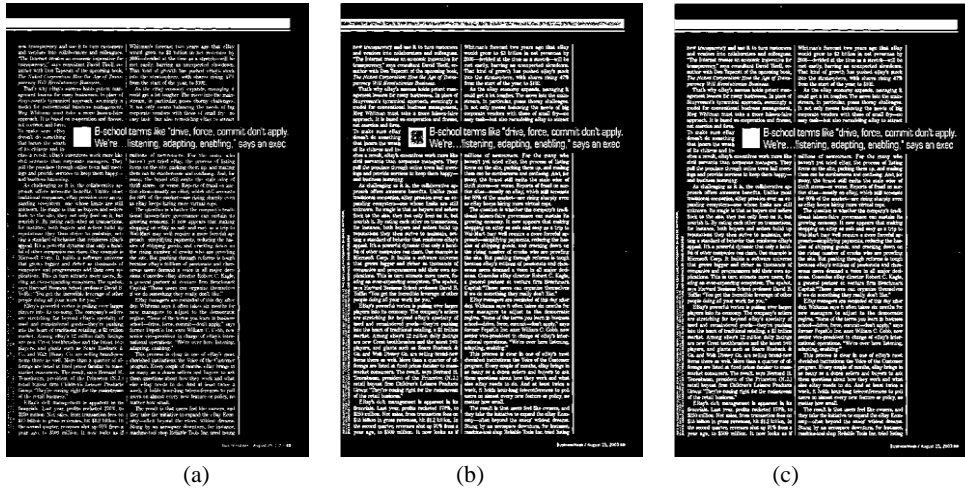
Fig. 1. Initialization input image (a) Otsu (b) Sauvola (c) Linear combining result between Otsu and Sauvola

The remainder of our paper consists of following sections. In Section 2, we introduce the fuzzy set theory. Section 3 presents the method for extracting component level features from connected component algorithm result. Section 4 provides detail for generating fuzzy energy matrix. Section 5 briefly demonstrates post-processing algorithm. In Section 6, we show the experimental results with public datasets, and we conclude our paper in Section 7.

## 2. FUZZY THEORY

A Fuzzy set is a set whose elements have levels of membership represented by a real number in the interval [0, 1]. In order to represent imprecise vectors, fuzzy set is very powerful tool that uses degree of variables rather than quantitative variables. For example weather condition may have a value such as '1: good' or '0: bad' that are not clearly defined because weather has various conditions. But the fuzzy set can represent binary weather conditions as meaningful classifications [1]. The definition of fuzzy set can be given as follows:

**Definition.** Let U be the universe of discourse, where $U = u_1, u_2, \ldots u_n$. A fuzzy set is determined depending on universe of discourse and sub intervals. The Fuzzy set $A$ can be defined on the universe of discourse follow as;

$$A = \frac{f_A(x,u_1)}{u_1} + \frac{f_A(x,u_2)}{u_2} + \cdots + \frac{f_A(x,u_n)}{u_n} \quad (1)$$

Where $f_A$ denotes the membership function of the fuzzy set A, $f_A: U \to [0, 1]$, and $f_A(x, u_i)$ denotes the degree of membership function of the crisp interval $u_i$ with input linguistic value x, and $f_A(x, u_i) \in [0,1]$ and $1 \leq i \leq n$, linguistic values x should be assigned to each fuzzy set [2]. Here $f_{A_i}(x, u_j)$ is a rectangular membership degree and it is shown in equation (2)

$$f_A(x, u_i) = \begin{cases} 1 & , x \in u_i \\ 0 & , otherwise \end{cases} \quad (2)$$

In equation the observations of features are converted to the fuzzy set has got the highest degrees of membership value. In our study, we changed quantitative features of text and non-text into degree of variables using Fuzzy set rules.



(a)



(b)

Fig. 2. Binarization with low-resolution image (a) original (b) binary result.

## 3. FEATURE EXTRACTION

The aim of this section is to extract features, which represent text and non-next. To extract features, an input image is down-sampled to 1/5 because the low-resolution document image is enough to segment document and can reduce time expense. We have utilized combined 2 binarization results using otsu [3] and sauvola [4]. Each binarization technique has advantages: the otsu can clearly keep figure region and the sauvola can keep table and separator therefore we used linear combination result with morphological filling operation as an input image. Fig. 1 illustrates this process. Given binary input image, the connected component algorithm is applied to binary image. In order to categorize between text and non-text, we proposed a feature extraction function (3) based on spatial properties of the connected components

$$F_i = \alpha H_i + \beta |H_i - W_i| \qquad (3)$$

Where *H,W,* height and width of component respectively. The $\alpha$ and $\beta$ are weight value. The function $F$ is designed to increase energy when component is non-text. Usually non-next has not uniformly variation in comparison with text and gap of scalar between $H$ and $W$ is bigger than text. The function F depends on $H$ value because our input image is low-resolution image so we often see distorted width features from touched small text horizontally whereas most of height features keep its property despite low-resolution environment as in Fig. 2. In this part, we just extract elemental features based on geometrical properties so the function F is not enough for categorizing component into text and non-text. The reason is because feature of several non-texts have similar features with big text. This problem is main challenge task in the document analysis field. We will figure out problem using Fuzzy Energy Matrix in the next section.

## 4. FUZZY ENERGY MATRIX (FEM)

In this section, the ultimate goal is to generate the Fuzzy Energy Matrix to classify text and non-text. The FEM is two-dimensional matrix representing the grade of membership of $F$ and location in the document. As mentioned in the previous section, the function $F$ is not enough to recognize between big text and figure so we included a location feature in the document. These combining features can provide that how many components, which have similar $F,$ are located in similar level space to us. The location of components is an efficient feature for recognizing big text from figure because most of big texts are arranged in a row. Therefore, we will generate the FEM using based on mentioned properties.

### 4.1 Fuzzification

After connected component algorithm, we extracts two features, which are height coordinate of vector and $F$, and then extracted quantitative features are transformed into "grade of membership" using universe of discourse and sub intervals. In the fuuzy set theory, this process is called "Fuzzification". The lengths of intervals used in partitioning universe of discourse are determined subjectively. The following Fuzzification algorithm can be given:

Step 1. Define the two universe of discourse and subintervals.
   Based on document height.

$$U = \left[ 1 , \frac{h}{2} \right], V = [ 1 , h ] \qquad (4)$$

Where h is document height. The $U$ and $V$ are universe of discourses based on function $F$ and image location, respectively.

Step 2. The length of intervals is defined as $n = 15$ and equation can be defined follows;

$$R_{bound_i} = k \left( \frac{Ud_{max}}{n} \right)$$

$$L_{bound_i} = (k - 1) \left( \frac{Ud_{max}}{n} \right), k, i=1....n \quad (5)$$

$$v_i, u_i = [ L\_bound_i , R\_bound_i ]$$

Where Ud represents each universe of discourse.
[ $R_{bound_i}$ ] and [ $L_{bound_i}$ ] are boundary (right/ left) of intervals.

Step 3. We fuzzified input features, and pseudocode can be given:

---

*Function Fuzzification*

*Input : 2 Main Features [F , Mh]    *Mh: location*
*       Number of component [N] ,*
*       Length of intervals [$v_i, u_i$ ,n]*
*Output : 2 Fuzzified vectors [Fv₁ ,Fv₂]*
 *for ( i =1 to n)*
   *{*

$$Fv_{1i} = \sum_{k=1}^{n} f_A(F_i, u_k)/u_k$$

$$Fv_{2i} = \sum_{k=1}^{n} f_A(Mh_i, v_k)/v_k$$

*}*
 *Return Fv₁ , Fv₂*

---

The fuzzified vectors ($Fv_1, Fv_2$) based on the length of intervals are applied to generate Fuzzy Energy Matrix in the next section.

### 4.2 Matrix Generation

The Fuzzy Energy Matrix represents 2-dimensional probability matrix for classifying between text and non-text based on Fuzzified vectors from previous section. In order to generate FEM, first we make a $n \times n$ matrix which represents fuzzy grade for *F, location*. And then whole of fuzzified vectors are assigned to level of matrix and then we count frequency corresponding to each grade. In the Table 1, we can observe that how many components are located in each fuzzy grade and assume that fuzzy level space including large energy is text. Whereas isolated small energy is non-text because when similar F values are located in similar grade in the Fuzzy Energy Matrix, it is strong feature representing the text. The integrated values in Table 2 are normalized to range 0~1 using sigmoid function [5]. The normalization equation can be given:

$$Energy_{xy} = \frac{1}{1+e^{-\gamma(matrix_{xy}-X)}} \qquad (6)$$

Where $\gamma$ is gradient of curve and $X$ is Threshold value considering text feature. In this study, we have used a X=4. If number of same F values are detected in same location more than X=4, it indicates a high probability with text. A pseudocode for FEM process can be given:

---

*Function Fuzzy Energy Matrix*
*Input : 2  Fuzzified vectors [Fv₁ ,Fv₂] ,Number of component [n] ,*
*        Length of intervals [$v_i, u_i$ ,n], Threshold X,*
*Output : Fuzzy Energy Matrix*

---

*Define FEM[n][n]; % Matrix for counting components*
*Define n_FEM[n][n] % Matrix for normalization using sigmoid*

*for(i =1 to n)*
  *{*
     *FEM[$Fv_{1i}$][ $Fv_{2i}$]= FEM[$Fv_{1i}$][ $Fv_{2i}$] + 1;*
  *}*

*for(i =1 to n)*
  *for(j =1 to n)*
    *{*
     *n_FEM[i][j]= $\frac{1}{1+e^{-\gamma(FEM[i][j]-X)}}$*
    *}*
*Return FEM, n_FEM*

In Fig. 3, we can see segmented non-text energy with green coordinates from the Table 1, Table 2. And the Fig. 3 (c) is a document energy corresponding to *n_FEM*. The text and non-text results are defined using document energy with T=0.5 (energy >0.5: text, energy <0 .5: non-text).


(a)               (b)


(c)


(d)               (e)

Fig. 3. Classification text using FEM (a) original (b) binary result (c) Energy map with 3D (d) text segmentation result (e) non-text segmentation result (green coordinates are fuzzified values * we can check energy in the Table 1, Table 2)

## 5. POST-PROCESSING

In this section, the main focus was introducing on five tasks which are Table detection, noise, separator, generating boundary for text.

### 5.1 Text boundary
In order to generate boundary region for text, we proposed a background extraction algorithm based on projection technique. This process is as follows:

Step 1. The morphology horizontal dilation algorithm is applied to the segmented text results (Fig. 4(a)).

Step 2. In the dilation result, projection technique with horizontal and vertical are utilized for extracting background. The stop condition is that when projection algorithm may be faced with text pixel. And then if length of projection result is more than (image height/2, width/2), the result is considered background. Fig. 4 shows estimated background.

Step 3. Finally, we observe non-text pixels in the black region Fig. 4 (c). If there are no non-text pixels, we made rectangle around text in the observing black region whereas the morphology result is just utilized for case of observed non-text.


(a)               (b)

(c)                                    (d)

Fig. 4. Text boundary generation (a) segmented text (b) morphology result (c) background extraction (d) final result

### 5.2 Heuristic filter for Non-text

First, in order to segment Table, our algorithm examines components, which has fuzzy grade $u_1$, are existed in the non-text region because most of table has a number of similar components arranged in the non-text region uniformly. In this study, when number of counted component $u_1$ in the non-text more then 10, it considers table. When table doesn't have lines, this way has strong advantage compare to traditional methods because we easily utilized geometric information of words arranged table using fuzzified values. Second, the separator is 2-type which are horizontal and vertical. A ratio between horizontal and vertical is utilized to detect separator. If ratio less than 0.2, we have considered the separator. Finally, the noise components are eliminated according to size of component.

Table 1. Fuzzy Energy Matrix with Document in Fig. 3

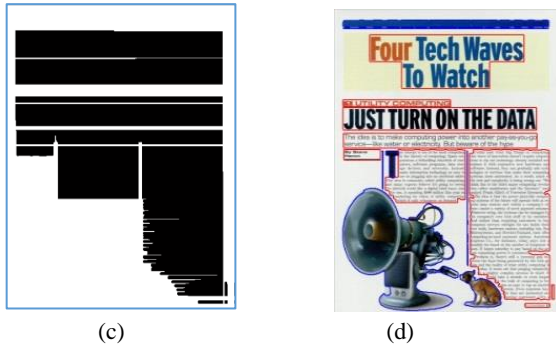|          | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ | $v_9$ | $v_{10}$ | $v_{11}$ | $v_{12}$ | $v_{13}$ | $v_{14}$ | $v_{15}$ |
|----------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|----------|----------|
| $u_1$    | 0     | 0     | 1     | 1     | 12    | 5     | 55    | 88    | 77    | 64       | 37       | 39       | 44       | 45       | 26       |
| $u_2$    | 0     | 4     | 7     | 4     | 0     | 11    | 0     | 1     | 0     | 0        | 0        | 0        | 0        | 0        | 1        |
| $u_3$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 1     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_4$    | 2     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_5$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_6$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 1        | 0        | 0        | 0        |
| $u_7$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 1        | 0        | 0        | 0        | 0        | 0        |
| $u_8$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_9$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{10}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{11}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{12}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{13}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{14}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{15}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |

Table 2. Normalization Fuzzy Energy Matrix using Sigmoid in Fig. 3

|          | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ | $v_9$ | $v_{10}$ | $v_{11}$ | $v_{12}$ | $v_{13}$ | $v_{14}$ | $v_{15}$ |
|----------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|----------|----------|
| $u_1$    | 0     | 0     | 0.1   | 0.1   | 1     | 1     | 1     | 1     | 1     | 1        | 1        | 1        | 1        | 1        | 1        |
| $u_2$    | 0     | 0.7   | 0.9   | 0.7   | 0     | 0.9   | 0     | 0.1   | 0     | 0        | 0        | 0        | 0        | 0        | 0.1      |
| $u_3$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0.1   | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_4$    | 0.3   | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_5$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_6$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0.1      | 0        | 0        | 0        |
| $u_7$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0.1      | 0        | 0        | 0        | 0        | 0        |
| $u_8$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_9$    | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{10}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{11}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{12}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{13}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{14}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $u_{15}$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |

## 6. EXPERIMENTAL RESULTS

We have implemented the proposed algorithm using MATLAB 2012 on an Intel(R) Core(TM) i5-4670 CPU 3.40GHz. The performance of the proposed method is evaluated on public ICDAR 2009 datasets because they contain 55 document images and have been utilized in most comparative works [12]. Also we have utilized our 100 document images, which are 50 English and 50 Korean for evaluation with ground truth mask.

The DICE system comprises two main steps. First each pixel is categorized primarily into machine-printed text. Second, isolated vectors are eliminated. Third, open and close operations are used to remove small regions. Finally, interior pixels are removed and contours of polygons are extracted [12]. The Fraunhofer system includes hybrid methods, which are

separator detection, page segmentation and Text line and region extraction [11]. The Tesseract is bottom-up method submitted by Ray Smith of Google Inc. This method includes binary morphology and connected component analysis, to estimate the type, which are text, image, separator and unknown. Two of the key methods employed include neighbourhood stroke-width measurement and appropriateness of overlap between adjacent connected components [13].

Fig. 5 demonstrates the performance using F-measure result with other state-of-the-art methods. As mentioned previously, our algorithm focused on processing time however our performance is also good as compare to traditional methods in the fig. 5, fig. 6 and average of proposed algorithm processing time is a 1.4 second in the MATLAB environment. Unfortunately, other state-of-the-art methods didn't mention their processing time. However the state-of-the-art methods are operated on the full-resolution image so we can guess that their processing time may well exceed 3 second because while implementing connected component and binarization algorithm in the full-resolution document, algorithm requires processing time of at least 2 second in the ICDAR 2009 datasets (PRIMA). Whereas proposed algorithm have used low-resolution document image (scale 1/5) to decrease time consuming and our performance is not so bad. Fig. 5 shows proposed algorithm performance with other state-of-the-art methods using F-measure.



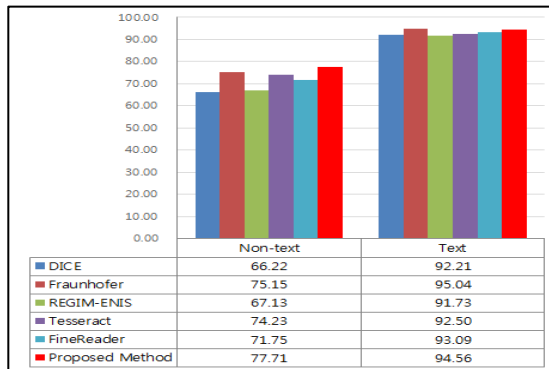| | Non-text | Text |
|---|---|---|
| DICE | 66.22 | 92.21 |
| Fraunhofer | 75.15 | 95.04 |
| REGIM-ENIS | 67.13 | 91.73 |
| Tesseract | 74.23 | 92.50 |
| FineReader | 71.75 | 93.09 |
| Proposed Method | 77.71 | 94.56 |

Fig. 5. F-measure result with ICDAR 2009 datasets

In the Fig. 6, the PRIMA measure for different region types for the four submitted methods. The PRIMA Lab (University of Salford, Manchester, United Kingdom) provides evaluation tool for PRMIA measure.
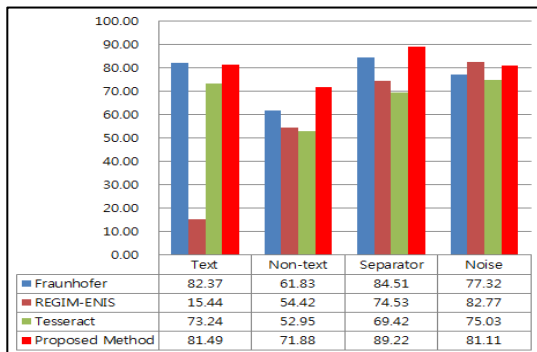


| | Text | Non-text | Separator | Noise |
|---|---|---|---|---|
| Fraunhofer | 82.37 | 61.83 | 84.51 | 77.32 |
| REGIM-ENIS | 15.44 | 54.42 | 74.53 | 82.77 |
| Tesseract | 73.24 | 52.95 | 69.42 | 75.03 |
| Proposed Method | 81.49 | 71.88 | 89.22 | 81.11 |

Fig. 6. PRIMA measure with ICDAR 2009 datasets



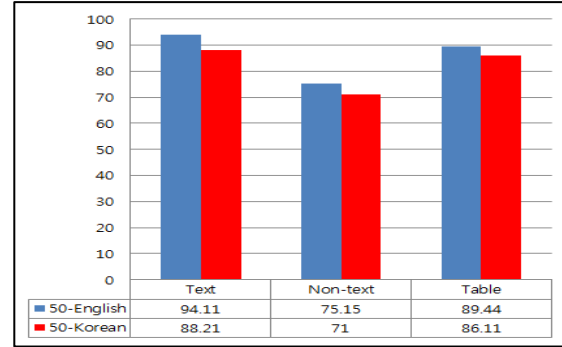| | Text | Non-text | Table |
|---|---|---|---|
| 50-English | 94.11 | 75.15 | 89.44 |
| 50-Korean | 88.21 | 71 | 86.11 |

Fig. 7. F-measure with 100 datasets with Korean and English

Fig. 8 illustrates visualization results of proposed method with 3D maps

## 7. CONCLUSION

We have presented a document image analysis method using FEM with low-resolution image. In particular, our algorithm performance and processing time are good compare to the state-of-art-methods however there are several drawbacks. First, the FEM algorithm can only detect text and non-text, in order to detect, table, separator, noise, it depends on heuristic post-processing algorithm. As mentioned previously, a heuristic algorithm is not best way. Therefore we will figure out the first drawback using improvement in $F$ equation. Second, proposed algorithm performance depends on initial binarization result. When character are touched each other in the binarization result, it may cause a critical effect on our proposed mechanism. In order to overcome second problem, we need extract a reasonable feature for recognizing touched character in the future works.

## ACKNOWLEDGEMENT

## REFERENCES

[1]  P. Bhupendra Kumar and S. Sanjay, "Integrated Fuzzy-HMM for project uncertainties in time-cost tradeoff problem," J. Applied Soft Computing, vol. 21, 2014, pp. 320-329.

[2]  S. M. Chen, "Forecasting enrollments based on fuzzy time series," Fuzzy Sets and Systems, vol. 81, 1996, pp. 311-319.

[3]  N. Otsu, "A Threshold Selection method from Gray Level Histogram," IEEE Transactions on Systems, Man and Cybernetics, vol. 9, no. 1, 1975, pp. 62-66.

[4]  J. Sauvola and M. Pietikainen, "Adaptive document image binarization," J. The Journal of the Pattern Recognition, 2000, pp. 225-236

[5] A. Fletcher and R. Kasturi, "A robust algorithm for text string separation from mixed text/graphics images," IEEE Trans. Pattern Anal. Mach. Intell., vol. 10, 1998, pp. 294-308.

[6] K. J. Anil and Y. Bin, "Document representation and its application to page decomposition," IEEE Trans. Pattern Anal. Mach. Intell., vol. 20, 1998, pp. 294-308.

[7] D. L. O'Gorman, "The document spectrum for page layout analysis," IEEE Trans. Pattern Anal. Mach. Intell., vol. 15, 1993, pp. 1162-1173.

[8] C. Laura, C. Ciro, and G. Przemyslaw, "Document page segmentation using neuro-fuzzy approach," Applied Soft Computing, vol. 8, 2008, pp. 118-126.

[9] K. Kise, A. Sato, and M. Iwata, "Segmentation of page images using the area Voronoi diagram," In Computer Vision Image Understanding, vol. 70, no. 3, 1998, pp. 370-382.

[10] D. H. S. Baird, "Backgroud structure in document images," Int. J. Pattern Recognition Artif. Intell., vol. 8, 1994, pp. 1013-1030.

[11] K. J. Anil and Y. Bin, "Document representation and its application to pagedecomposition," IEEE Trans. Pattern Anal. Mach. Intell., vol. 20, 1998, pp. 294-308.

[12] A. Antonacopoulos, S. Pletschacher, D. Bridson, and C. Papadopoulos, "ICDAR2009 Page Segmentation Competition", Proc. ICDAR, 2009, pp. 1370-1374.

[13] Ray Smith, "Hybrid Page Layout Analysis via Tab-Stop Detection", Proc. ICDAR, Barcelona, Spain, 2009, pp. 241-245.

**Kang Han Oh**
He received the B.S in Computer Science from Honam University in 2010. And he received the M.S in Electronic & Computer Engineering at Chonnam National University in 2013. He has been taking the Ph.D course in Electronics & Computer Engineering at Chonnam National University, Korea. His research interests are pattern recognition, machine learning and Image processing.

**Soo Hyung Kim**
He received his B.S degree in Computer Engineering from Seoul National University in 1986, and his M.S and Ph.D degrees in Computer Science from Korea Advanced Institute of Science and Technology in 1988 and 1993 respectively. From 1990 to 1996, he was a senior member of research staff in Multimedia Research Center of Samsung Electronics Co., Korea. Since 1997, he has been a professor in the Department of Computer Science, Chonnam National University, Korea. His research interests are pattern recognition, document image processing, medical image processing, and ubiquitous computing.
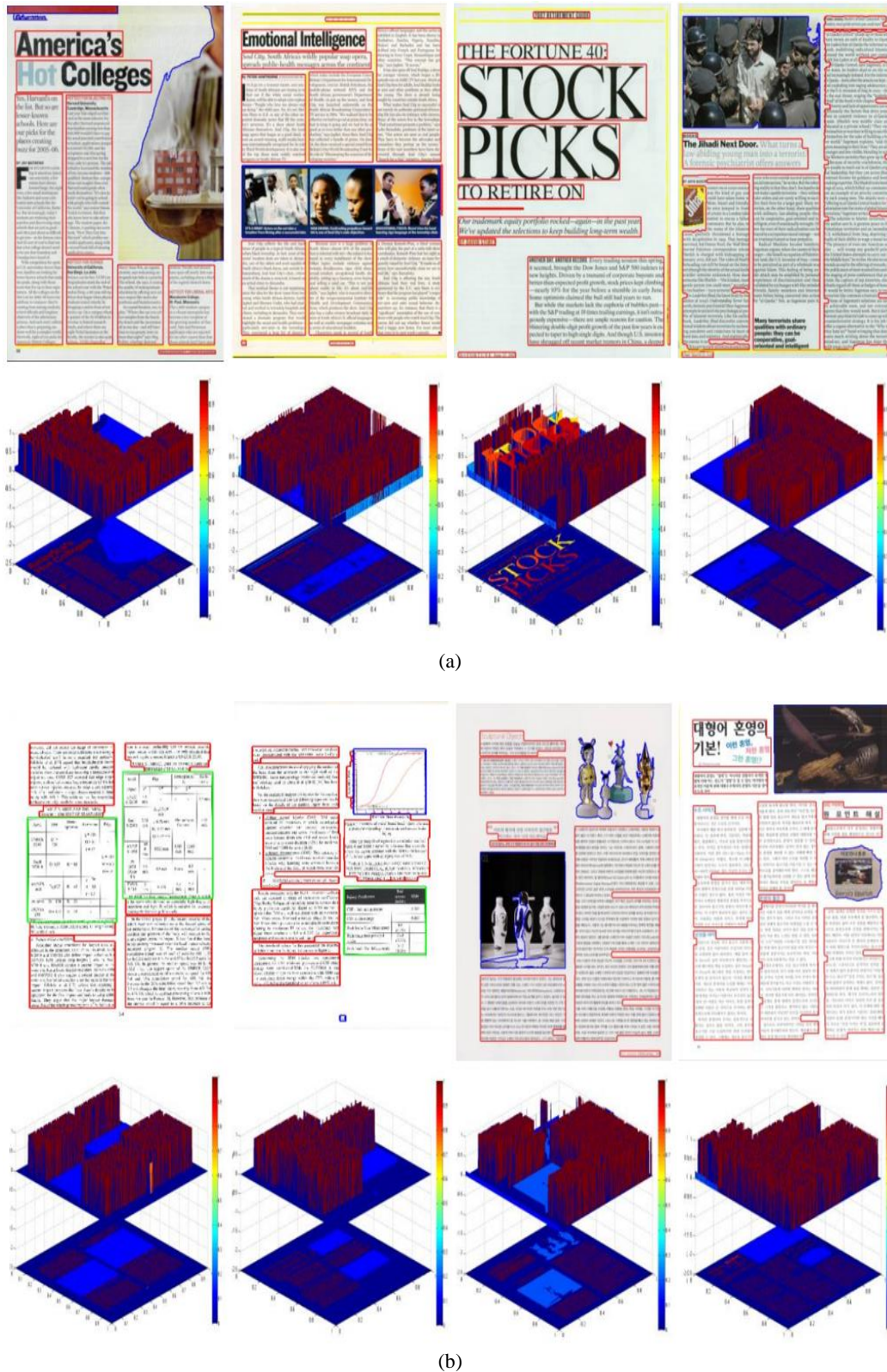
(a)



(b)

Fig. 8. Visualization results of proposed method with 3D map (a) ICDAR 2009 (b) Our datasets.