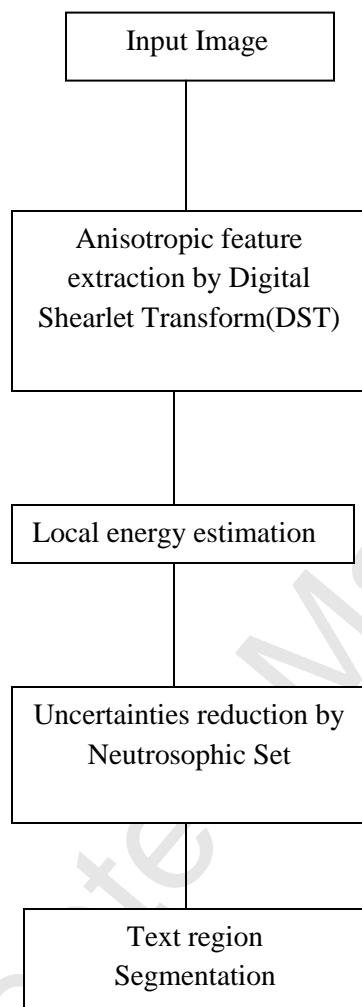


Highlights

- In this paper, a novel text region segmentation method based on Digital Shearlet Transform(DST) is proposed.
- The uncertainties in the DST features are handled by Neutrosophic set theoretic approach.
- The proposed method shows robustness under different perturbations.

Accepted Manuscript

Graphical Abstract

Accurate segmentation of complex document image using Digital Shearlet Transform with Neutrosophic Set as uncertainty handling tool

Soumyadip Dhar^a, Malay K. Kundu^b

^a*RCC Institute of Information Technology ,Kolkata-700015,India, rccsoumya@gmail.com*

^b*Indian Statistical Institute,Kolkata-700108,India, malay@isical.ac.in*

Abstract

In any image segmentation problem, there exist uncertainties. These uncertainties occur from gray level and spatial ambiguities in an image. As a result, accurate segmentation of text regions from non-text regions (graphics/images) in mixed and complex documents is a fairly difficult problem. In this paper, we propose a novel text region segmentation method based on Digital Shearlet transform(DST). The method is capable of handling the uncertainties arising in the segmentation process. To capture the anisotropic features of the text regions, the proposed method uses the DST coefficients as input features to a segmentation process block. This block is designed using the Neutrosophic set(NS) for management of the uncertainty in the process. The proposed method is experimentally verified extensively and the performance is compared with that of some state-of-the-art techniques both quantitatively and qualitatively using benchmark dataset.

Keywords: Shearlet , Digital Shearlet transform(DST), neutrosophic set, uncertainty handling, document image, segmentation

1. Introduction

Due to the rapid development of digital technology over the last few decades, usage of digital documents becomes a common practice. The practice is due to low cost for storage, easy storage, and transportability of the digital documents. Automation of digital document analysis involves the segmentation of text regions from non-text regions/graphics.

Proper text and non-text region segmentation help to store and retrieve the text documents

Preprint submitted to Applied Soft Computing

August 11, 2017

16 efficiently. In the area of text and graphics region segmentation, a lot of research work has
17 already been reported so far. These works include segmentation of structurally complex
18 (color and gray level) document and the camera captured natural scene images.

19 The conventional approaches to image segmentation consist of hard partitioning the
20 image space into meaningful regions by extracting its features. But, the regions are not
21 crispy defined due to incomplete or imprecise input information, the ambiguity/vagueness
22 in an input image, ill-defined and /or overlapping boundaries among the regions, and the
23 indefiniteness in defining/extracting features. An image possesses ambiguity within each
24 pixel because of the possible multi-valued levels of brightness. The uncertainty in an image
25 pattern may be explained in terms of gray level [1] or spatial ambiguity [2] or both. Gray
26 level ambiguity means "indefiniteness" in defining a pixel as white or black(i.e object or
27 background) and it considers global information. Spatial ambiguity refers to "indefiniteness"
28 in shape and geometry(e.g defining centroid, sharp edge etc.) within the image[3]and it takes
29 care of local information. The uncertainties may increase in images due to noise, rotation
30 and dynamic gray level change. Due to the uncertainties, the proper feature generation for
31 segmentation process is difficult. In order to locate the segmentation boundary between
32 the regions accurately, it is necessary that the segmentation process should be capable of
33 handling the uncertainties in an effective manner. The ultimate output of the process will be
34 associated with least uncertainty and the output retains as much of the 'information content'
35 of the image as possible. Thus, it is natural and convenient to avoid committing ourselves
36 to a specific(hard) decision about the segments. The segments should be represented by the
37 subsets which are characterized by the degree to which each pixel belongs to them.

38 *1.1. Related work*

39 Conventional methods for text region segmentation include top-down and bottom-up ap-
40 proaches. In the top-down approach, at first, the probable text regions are identified. Then
41 the regions are split into paragraphs, text-line, and words. X, Y-cut algorithm [4], which
42 works by detecting white space using horizontal and vertical projection, is one of the most
43 popular top-down approaches. Chen et al. [5] proposed a top-down approach using canny

44 filter and a vertical, horizontal edge map to detect the text regions. Alex Chen et al. [6] in
45 his top-down approach identified the text regions of an image by their statistical properties.
46 On the contrary, bottom-up approach attempts to cluster the pixels and then group together
47 to obtain the text and non-text regions. A bottom-up approach by a block-wise text pixel
48 segmentation method was proposed by Haneda et al. [7]. Here the initial segmentation was
49 refined by a connected component classification (CCC). Yi et al. [8] developed a bottom-up
50 algorithm for segmentation of text in natural scenes based on adjacent character grouping
51 method and text line grouping method. Their method suffered from difficulties like the seg-
52 mentation failure for a ligature, multicolored text regions, and text sets comprise of less than
53 three characters. Currently, the bottom-up approach, based on Maximally Stable External
54 Regions (MSERs) [9] has gained popularity. The MSERs are based on the idea of taking re-
55 gions which stay nearly same through a wide range of thresholds. Gomez et al. [10] proposed
56 an algorithm for text extraction based on MSERs and human perception of text in natural
57 scenes. MSERs based text region segmentation was also used by [11, 12, 13, 14]. Another
58 bottom-up approach based on seed generation was proposed by Bai et al. [15]. Cho et al. [16]
59 proposed a text region segmentation based on MSERs and canny edge detector. However,
60 the methods described above, are sensitive to character size, scanning resolution, inter-line,
61 and inter-character spacing. They suffer from low accuracy rate when prior knowledge about
62 the image content is not available.

63 Apart from the above mentioned methods, some researchers used deep learning based
64 methods for the text region detection where Convolutional Neural Network(CNN) was used
65 for text region classification. He et al. [17] used CNN for text component filtering and
66 incorporated Contrast-Enhanced MSERs(CE-MSERs) for text region detection. Huang et
67 al. [18] proposed a deep learning based method, which integrated the MSERs and CNN. In
68 the method, MSERs was used for text region detection and CNN based classifier was utilized
69 for true identification of text regions. However, the method and other deep learning based
70 methods were computationally costly, as they required tremendous data for training and
71 the methods were text dependent. The performance boost of deep learning methods may
72 be due to the training with huge amount of data which may not be publicly available [19].

73 On the other hand they could only detect horizontal or near horizontal text regions. These
74 two limitations may restrict the practical application of deep learning based methods [19].

75 Another approach is used by the researchers where it is assumed that text regions have
76 distinctive texture properties. The text region texture properties are quite different from
77 the background texture properties. So segmentation of text from non-text regions is treated
78 as a texture segmentation problem. Zhou et al. [20] combined the three different texture
79 features, such as oriented gradient (HOG), mean of gradients (MG) and local binary pat-
80 terns(LBP) for text region segmentation. Some of the techniques of this approach use Gabor
81 filter, Wavelet transforms etc., as a representation tool for feature extraction. The proposed
82 method uses a similar methodology for text and non-text region segmentation. The dyadic
83 wavelet based text region segmentation was reported in [21, 22]. But the dyadic wavelet
84 coefficients do not yield rotation invariant features. The limitation can be overcome using
85 Gabor filter. Chan et al. [23] proposed a method for text region segmentation that involves
86 computation of local energy for texture using a bank of orthogonal pairs of Gabor filter.
87 Gabor filter bank was also used for generating features by Nirmal et al. [24] to detect the
88 text regions. The limitation of the methods is that Gabor filter banks come with the high
89 computational costs. To overcome the limitation of both the dyadic wavelet transform and
90 Gabor filter, M-band wavelet transform and M-band wavelet packet transform are preferred
91 by many researchers. Acharyya et al. [25] used M-band wavelet transform to segment the
92 text regions from gray images. Kumar et al. [26] designed matched wavelet and Markov
93 Random Field model for document image segmentation. But, the Wavelet transform can
94 only handle the point singularities and cannot capture the curve like features properly.

95 The methods described above and most of the methods reported in the literature do
96 not have the capacity to capture the anisotropic features (edges, curves) of the text regions
97 properly. Moreover, they cannot handle the uncertainties inherent in an image and the
98 uncertainties that can arise in feature generation [27]. Maji. et al. [28] used a rough-fuzzy
99 clustering technique on M-band packet wavelet features to handle the uncertainties. Kundu
100 et al. [29] proposed a text region segmentation method combining M-band packet wavelets
101 and fuzzy c-means algorithm. But their uncertainty handling models were used during the

102 segmentation only and had no effect on feature generation. As a result, the performance of
103 the methods did not achieve the expected level of accuracy. So, that is the major motivation
104 of the current investigation for a better solution.

105

106 1.2. Proposed Method

107 We propose an accurate text region segmentation method which can capture the curve
108 and edge like feature of text regions in an image and also handle the uncertainties in the
109 image during feature generation. The features help to represent the two different textures
110 of text and non-text regions properly. The method is based on Neutrosophic Set(NS) and
111 multi-resolution analysis of the image using Digital Shearlet Transform(DST) transform [30].
112 The motivation for using DST is that, it can handle the anisotropic features better than the
113 wavelet transform. Moreover, it has better multi-directional, shift invariant and excellent
114 multiscale image decomposition property than the curvelet transform and contourlet trans-
115 form [31, 30].

116 Neutrosophic Set was proposed by Florentin Smarandache as a new branch of philosophy
117 dealing with the origin, nature and scope of neutralities [32]. It has a powerful capacity to
118 deal with the uncertainty and is better than other uncertainty handling model [33]. The
119 concept was successfully used in image thresholding [34], image denoising [35], image seg-
120 mentation [33] and color texture image segmentation [36]. In NS theory, every event has
121 not only a certain degree of truth but also a falsity degree and an indeterminacy degree that
122 have to be considered independently from each other [32] and represented as true, false and
123 indeterminate set respectively.

124 In the proposed feature extraction method we use the DST to transform the image into
125 different shearlets translates(sub-bands)which have different scales and orientations. For
126 each shearlet coefficient in a sub-band, the energy is computed from the transform coeffi-
127 cients over a overlapping window (size adaptively varied) around each pixel. The transform
128 coefficients are itself rotational invariant. The anisotropic features of the text regions are
129 captured by the DST efficiently. The anisotropic features represent the text regions irrespec-

130 tive of the character size, inter-line, and inter-character spacing. Thus, the DST captures
131 the texture property of the text regions to differentiate it from the non-text/graphics re-
132 gions. The texture represents the assembly of text fonts, unlike the conventional methods
133 which used textual and graphical attributes like font size, text line orientation etc for text
134 region segmentation. So the proposed method can work in a generic environment. But,
135 DST may introduce uncertainties in features due to the presence of redundant feature (mul-
136 tiple numbers of shearlet sub-bands) information. This is in addition to conventional type
137 uncertainties present due to the discrete gray level ambiguity, the spatial ambiguity which
138 may increase due to different perturbations and low scanning resolutions etc. These are
139 major reasons for inaccuracy for any segmentation process in locating true segmentation
140 boundaries. So in order to tackle this problem effectively, we have to use the NS as an
141 uncertainty handling model, which can reduce error due to uncertainties for achieving the
142 better segmentation.

143 In the proposed uncertainty handling scheme, at the beginning, the features(local energy
144 of shearlet coefficients) are mapped into NS domain in order to classify them into three
145 different sub sets, the true, indeterminate and false. In this stage, the uncertainties are
146 reduced iteratively with the increase of values of entries in true and false subsets along
147 with the reduction values of entries in indeterminate sub sets. After this reduction of the
148 uncertainties, the feature selection is done on the true subsets by an unsupervised feature
149 selection algorithm to generate finer features. This step helps to reduce the uncertainties
150 due to the presence of a large number of redundant features. With this re-categorized fea-
151 ture subset, the final segmentation is done using NS based clustering. This is why in the
152 proposed method the two stage uncertainty handling capabilities is expected to give better
153 results than the other existing uncertainty handling models having no such provision of suc-
154 cessive uncertainty reduction mechanism.

155

156 *1.3. Novelty of the proposed method*

157 The novelties or the major contributions of the of the paper are

- 158 (1) : We propose an efficient texture based method by the DST for text region segmen-
159 tation. The highly efficient rotation and translation invariant anisotropic features
160 generated by the shearlets in DST help to capture the different textures of text and
161 non-text regions in a complex background. Thus, it greatly improves the performance
162 of text region segmentation over the state-of-the-art methods.
- 163 (2) :In the proposed method, the uncertainties in the sub-bands of the DST are handled
164 by the neutrosophic set. The uncertainties are due to spatial and gray level ambigu-
165 ity in an image. Moreover, additional uncertainties are introduced due to redundant
166 information generation of the shearlet sub-bands in different scales. The uncertainties
167 in the NS domain are handled in two steps. In the first step, the uncertainties in a
168 feature itself are reduced and in the second step uncertainties during the segmentation
169 are handled. For this in the first step, iteratively the uncertainties in each feature
170 are reduced in NS domain. On the top of that to generate the finer features, features
171 selection in the NS domain based on Maximal information compression index(MCI) is
172 done. In the second step, uncertainty during the segmentation is reduced by cluster-
173 ing in the NS domain. This two-step uncertainty reduction process in text/non-text
174 segmentation is more powerful than the conventional methods which handle the uncer-
175 tainties by fuzzy-c-means or rough-fuzzy c-means with no such provision of two-step
176 uncertainty reduction. They handle the uncertainties only during the segmentation
177 process. Such process has no effect on feature generations.
- 178 (3) : The proposed method is tested under different perturbations i.e noise corruption,
179 rotation, and dynamic gray level changes. Compared to the state-of-the-art methods,
180 our method shows satisfactory robustness under the different perturbations.

181 The paper is organized as follows. Section 2 describes DST, representation of NS components
182 and unsupervised feature selection. The proposed method and the algorithm are presented
183 in the Section 3. The Section 4 comprises of results discussion, comparison with other
184 methods and performance evaluation.

185 **2. Theoretical Preliminaries**

186 *2.1. Shearlet system*

187 Shearlet systems are designed to efficiently encode anisotropic features such as singular-
 188 ities concentrated on lower dimensional embedded manifolds. To achieve optimal sparsity,
 189 shearlets are scaled according to a parabolic scaling law encoded in the parabolic scaling
 190 matrix A_a , $a > 0$ and exhibit directionality by parameterizing slope encoded in the shear
 191 matrix S_s , $s \in \mathbb{R}$, defined by

192

$$A_a = \begin{bmatrix} a & 0 \\ 0 & \sqrt{a} \end{bmatrix} \text{ and } S_s = \begin{bmatrix} 1 & s \\ 0 & 1 \end{bmatrix}$$

193 For appropriate choices of the shearlet $\psi \in L^2(\mathbb{R}^2)$, the Continuous Shearlet Transform

$$\mathcal{SH}_\psi : f \rightarrow \mathcal{SH}_\psi f(a, s, t) = \langle f, \psi_{ast} \rangle \quad (1)$$

194 is a linear isometry from $L^2(\mathbb{R}^2)$ to $L^2(\mathbb{S})$. Hence, shearlet systems are based on three
 195 parameters: $a > 0$ being the scale parameter measuring the resolution level, $s \in \mathbb{R}$ being the
 196 shear parameter measuring the directionality, and $t \in \mathbb{R}^2$ being the translation parameter
 197 measuring the position. When $s < |1|$, this produces the cone adapted Continuous Shearlet
 198 Transform. It allows an equal treatment of all directions in contrast to a slightly biased
 199 treatment by the Continuous Shearlet Transform.

200 A discrete shearlet transform for $\psi \in L^2(\mathbb{R}^2)$, is a collection of functions of the form

$$\psi_{j,k,m} = 2^{3j/4} \phi(S_k A_{2^j} \cdot -m) : j \in \mathbb{Z}, k \in K \subset \mathbb{Z}, m \in \mathbb{Z}^2 \quad (2)$$

201 where K is a carefully chosen indexing set of shears. Note that the shearing matrix S_k maps
 202 the digital grid \mathbb{Z}^2 onto itself, which is the key idea for deriving a unified treatment of the
 203 continuum and digital setting. The discrete shearlet system defines a collection of waveforms
 204 at various scales j , orientations controlled by k , and locations dependent on m . To avoid the
 205 biased treatment of directions which the discrete system inherit, the cone adapted shearlet
 206 system is defined as

$$\mathcal{SH}(\phi, \psi, \tilde{\psi}) = \{\phi(\cdot - m) : m \in \mathbb{Z}^2\} \cup \{\psi_{j,k,m}, \tilde{\psi}_{j,k,m} : j \geq 0, |k| \leq \lceil 2^{j/2} \rceil, m \in \mathbb{Z}^2\} \quad (3)$$

207 where $\tilde{\psi}_{j,k,m}$ is generated from $\psi_{j,k,m}$ by interchanging both variables, and $\psi_{j,k,m}$, $\tilde{\psi}_{j,k,m}$ and
 208 ϕ are L^2 functions.

209 2.2. Digital Shearlet Transform

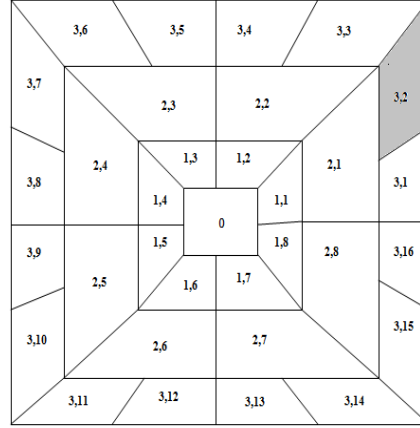


Figure 1: Digital Shearlet transform of an image into 4 levels. The 0th level is the low frequency shearlet. Each index represents the level and orientation of the shearlet in that level e.g the shearlet (3,2) represents the level 3 with orientation 2.

210 In the proposed method the Digital Shearlet transform [37] is based on cone adapted dis-
 211 crete shearlet system with compactly supported shearlets. In compactly supported shearlet
 212 system generator, it is conjectured that no tight shearlet frame exist. The DST is con-
 213 structed in three steps: (1)A non-separable structure of shearlet generation, (2)digitization
 214 of shearlet operators and (3)generation of digital shearlet filter.

215 A non-separable structure of shearlet generator $\hat{\psi}$ is defined as

$$\hat{\psi}(\xi) = P\left(\frac{\xi_1}{2}, \xi_2\right)\psi_1 \otimes \phi_1(\xi) \quad (4)$$

216 where P is a 2d directional filter, ϕ_1 is 1D scaling function associated with wavelet MRA
 217 and ψ_1 is the corresponding 1D wavelet function and $\xi = (\xi_1, \xi_2) \in \mathbb{R}^2$, $\hat{\mathbb{R}}^2$ being the plane
 218 in the Fourier domain. Because of $\psi_{j,k,m} = \psi_{j,0,m}(S_{k/2^{j/2}})$, two ingredients are required for
 219 digitization of shearlets, Digital shearlet filter $\psi_{j,0}^d$ and digital shearlet operator $S_{k/2^{j/2}}^d$. For
 220 any discrete 2D signal the shearlet operator is defined as $S_{k/2^{j/2}}^d(x) = ((x_{\uparrow 2^{d\alpha_j}} * h_{j/2})(S_k \cdot) * 1$

221 $\overline{h_{j/2}} \downarrow 2^{d_{\alpha_j}}$, where $x \in L^2(\mathbb{Z}^2)$, $\uparrow 2^{d_{\alpha_j}}$, $\downarrow 2^{d_{\alpha_j}}$ and $*_1$ are upsampling, downsampling and
 222 convolution operator along X axis and $d_{\alpha_j} = \lceil j(2 - \alpha_j)/2 \rceil$. Here the $\alpha_j (= 1)$ measures the
 223 degree of anisotropy for each scale $j \geq 0$. Then the digital shearlet filter is given by

$$\psi_{j,k}^d = S_{k/2^{j/2}}^d(x) * p_j * (g_{J-j} \otimes h_{J-j/2}) \quad (5)$$

224 Where p_j s are the Fourier coefficients of $P(2^{J-j-1}\xi_1, 2^{J-\frac{j}{2}}\xi_2)$, $h_{J-j/2}$ is the low pass filter
 225 associated with the scaling function ϕ_1 , g_{J-j} is the corresponding high pass filter, associated
 226 with the wavelet function ψ_1 and $J \in \mathbb{N}$ is the highest scale to be considered (i.e $j < J$ for all
 227 shearlets $\psi_{j,k,m}$). Now the digital shearlet transform of a digital signal $f \in L^2(\mathbb{Z}^2)$ is given
 228 by

$$\text{DST}_{j,k,m}^{2D} = \mathcal{SH}_{j,k,m}^d = (\overline{\psi_{j,k}^d} * f) \quad (6)$$

229 for $j \in \{0, J-1\}$ and $|k| < \lceil 2^{j/2} \rceil$

230 The shearlet systems use regular translations on the integer lattice with the shearing
 231 operations. The system provides the structure of the integer grid with directionality. Ob-
 232 viously, these two properties lead to an implementation of the digital shearlet transform
 233 exploiting discrete convolutions. This allows a shift-invariant transform by simply skipping
 234 the anisotropic downsampling. As a result, digital shearlet transform is highly redundant.
 235 The shearlet sub-bands generated by the DST have the same size of the input. Mainly, in
 236 two ways the uncertainties (Figure 2) are to be handled in the features generated by the
 237 DST. They are

- 238 1 : Each shearlet coefficient (feature) of a shearlet sub-band represents one pixel in one
 239 scale and orientation. Thus, the contrast between the coefficients in each shearlet
 240 sub-band should be increased to reduce the uncertainty due to gray and spatial level
 241 ambiguity in an image.
- 242 2 : The total number of shearlet sub-band generated in DST is $\sum_{j=0}^{n_{scales}-1} 2^{\lceil j/2 \rceil + 2}$, where
 243 n_{scales} represents the number of scales of the DST. The redundant information due
 244 to different scales and directions also increases the uncertainty. Proper selection of the
 245 sub-bands reduces the uncertainty.

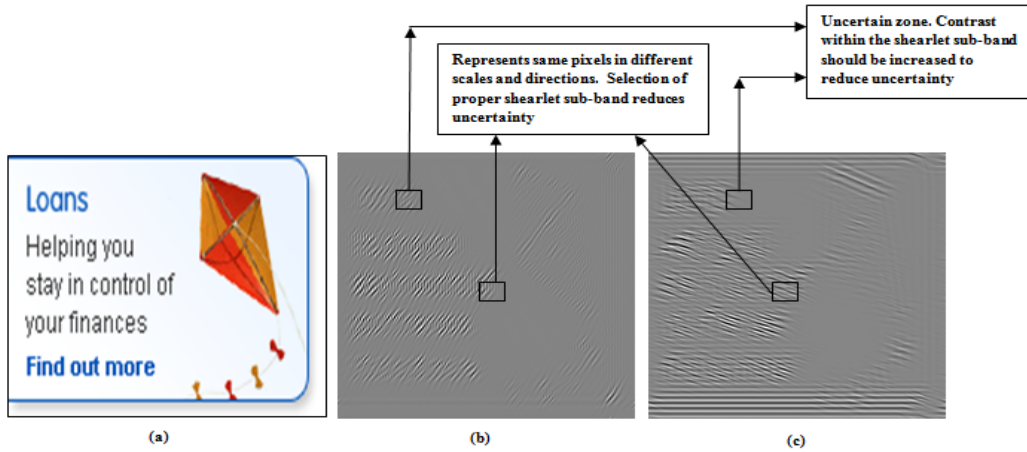


Figure 2: (a) Original image of size 256×256 (b) and (c) shows the two shearlets (size 256×256) by the DST of the image in two different resolutions and orientations with uncertain zones.

246 Apart from the above, uncertainties may occur during segmentation of the features. These
 247 uncertainties make it difficult to decide the correct class of a pixel. The uncertainties also
 248 increase due to different perturbations. To handle the above uncertainties the Neutrosophic
 249 set and Neutrosophic logic are used in the proposed method.

250 2.3. Neutrosophic set (NS):

251 Let U be a universe of discourse, and a neutrosophic set $A \subset U$. An element x from U
 252 is noted with respect to A as $x(T, I, F)$ and belongs to A in the following way [32]:
 253 x 's degree of belongings in the true subset T is $t\%$, in the indeterminate subset I is $i\%$ and
 254 in the false subset F is $f\%$. T, I and F are real standard or non-standard subsets in the
 255 open interval of $]^{-0}, 1^{+}[$
 256 with $supT = t_{sup}, infT = t_{inf}, supI = i_{sup}, infI = i_{inf}, supF = f_{sup}, infF =$
 257 f_{inf} and $n_{sup} = t_{sup} + i_{sup} + f_{sup}, n_{inf} = t_{inf} + i_{inf} + f_{inf}$. In the Neutrosophic
 258 Logic [32], the truth (T) and the falsity (F) and the indeterminacy (I) can be any numbers
 259 in $[0, 1]$, then $0 \leq T + I + F \leq 3$. For most real world applications, $T, I, F \subset [0, 1]; t + f = 1$
 260 and $i \in [0, 1]$ [33]. In the next section, we will discuss how neutrosophic set can be used as
 261 an uncertainty handling tool in the context of image segmentation problem.

262 2.4. Neurotrophic image representation for handling uncertainty :

Let $r(i, j)$ be the gray level of the (i, j) th pixel of a $P \times Q$ dimensional L level image $B = [r(i, j)], i = 1, 2, \dots, P, j = 1, 2, \dots, Q$. When the image is mapped into the neutrosophic domain, the corresponding image is called a neutrosophic image B_{NS} [36]. It is characterized by three subsets T, I and F . A pixel in an NS domain can be represented as (t, i, f) where the pixel is t true, i indeterminate and f false, and t, i, f belongs to the true subset T , the indeterminate subset I and the false subset F respectively. A pixel $B(i, j)$ in the image domain is mapped into NS domain as

$$B_{NS}(i, j) = \{T(i, j), I(i, j), F(i, j)\}$$

263 where $T(i, j)$, $I(i, j)$, and $F(i, j)$ are the membership values t, i, f belonging to true subset
 264 T , indeterminate subset I and false subset F respectively. They are represented by the
 265 following equations.

$$T(i, j) = 1 - \frac{\bar{r}_{max} - \bar{r}(i, j)}{\bar{r}_{max} - \bar{r}_{min}} \quad (7)$$

$$I(i, j) = 1 - \frac{d_{max} - d(i, j)}{d_{max} - d_{min}} \quad (8)$$

$$d(i, j) = \text{abs}(r(i, j) - \bar{r}(i, j)) \quad (9)$$

$$d_{max} = \max\{d(i, j) | i \in P, j \in Q\} \quad (10)$$

$$d_{min} = \min\{d(i, j) | i \in P, j \in Q\} \quad (11)$$

$$F(i, j) = 1 - T(i, j) \quad (12)$$

271 where $\bar{r}(i, j)$ is the mean value of the pixels within a local overlapping window of size
 272 $w \times w$ around $r(i, j)$ and its position is at the center of the local window. \bar{r}_{max} and \bar{r}_{min}
 273 are the maximum and minimum values of $\bar{r}(i, j) \forall i \in P, \forall j \in Q$. $d(i, j)$ is the absolute
 274 value of the difference between element $r(i, j)$ and its local mean value $\bar{r}(i, j)$. The value
 275 of $I(i, j)$ is employed to measure the indeterminacy degree of a pixel $B_{NS}(i, j)$. Here $I(i, j)$
 276 gives the degree of uncertainties of deciding brightness of the pixel. The indeterminacy of
 277 the NS image, which is measured by the entropy of the indeterminate subset, represents
 278 the uncertainties present in the gray level image. The changes in T and F influence the

279 distribution of element in I and the entropy of I . The NS domain representation of a gray-
 280 scale image can be found in Appendix A.

281 In NS domain, two operations, α – mean and β – enhancement, are used to reduce the
 282 indeterminacy of the neutrosophic image. After these two operations, the image becomes
 283 more uniform and homogeneous, and more suitable for segmentation. In the next two
 284 subsections, we will discuss the α – mean and β – enhancement operations mathematically.

285 2.5. The α – mean operation :

286 The α – mean operation [33] transforms neutrosophic image pixel $B_{NS}(i, j)$ to $B_{NS\alpha}(i, j) =$
 287 $\{\bar{T}(i, j), \bar{I}(i, j), \bar{F}(i, j)\}$ where \bar{T}, \bar{I} and \bar{F} are the true, indeterminate and false subsets of
 288 $B_{NS\alpha}$. In this operation, it is checked whether the indeterminate membership value of a pixel
 289 is higher than a predefined value α where $0 \leq \alpha \leq 1$. If so, the value is reduced by making
 290 the pixel homogeneous with neighboring pixels. Thus, the α – mean operation reduces the
 291 uncertainty due to the spatial ambiguity. It is represented as

$$\bar{T}(i, j) = \begin{cases} T(i, j) & \text{if } I(i, j) < \alpha \\ \bar{T}_\alpha(i, j) & \text{if } I(i, j) \geq \alpha \end{cases} \quad (13)$$

292 where $\bar{T}_\alpha(i, j)$ is the mean value of the pixels within a local overlapping window of size $w \times w$
 293 around $T(i, j)$ and its position is at the center of the local window. Similar operations are
 294 done for $\bar{F}(i, j)$. After the operation on the subset T, the indeterminate subset becomes

$$\bar{I}(i, j) = 1 - \frac{\bar{d}_{Tmax} - \bar{d}_T(i, j)}{\bar{d}_{Tmax} - \bar{d}_{Tmin}} \quad (14)$$

$$\bar{d}_T(i, j) = \text{abs}(\bar{T}(i, j) - \bar{\bar{T}}(i, j)) \quad (15)$$

$$\bar{d}_{Tmin} = \min\{\bar{d}_T(i, j) | i \in P, j \in Q\} \quad (16)$$

$$\bar{d}_{Tmax} = \max\{\bar{d}_T(i, j) | i \in P, j \in Q\} \quad (17)$$

298 where $\bar{d}_T(i, j)$ is the absolute value of the difference between the $\bar{T}(i, j)$ and its local mean
 299 value $\bar{\bar{T}}(i, j)$. $\bar{\bar{T}}(i, j)$ is calculated over a local overlapping window of size $w \times w$ around each
 300 $\bar{T}(i, j)$ after the α – mean operation on NS image.

301 2.6. The β – enhancement operation

302 In NS domain, the β –enhancement operation [33] is used to enhance the true membership
 303 value $T(i, j)$. The value is enhanced if its corresponding $I(i, j)$ value is greater than a
 304 predefined value β where $0 \leq \beta \leq 1$. The operation reduces the uncertainty due to gray
 305 level ambiguity by lowering the indeterminacy in the neutrosophic image. The β -enhanced
 306 image $B_{NS\beta}$ is defined as

$$B_{NS\beta}(i, j) = (T'(i, j), I'(i, j), F'(i, j)) \quad (18)$$

$$T'(i, j) = \begin{cases} T(i, j) & \text{if } I(i, j) < \beta \\ T'_\beta(i, j) & \text{if } I(i, j) \geq \beta \end{cases} \quad (19)$$

$$T'_\beta(i, j) = \begin{cases} 2T^2(i, j) & \text{if } T(i, j) < 0.5 \\ 1 - 2(1 - T(i, j))^2 & \text{if } T(i, j) \geq 0.5 \end{cases} \quad (20)$$

$$I'(i, j) = 1 - \frac{d'_{Tmax} - d'_T(i, j)}{d'_{Tmax} - d'_{Tmin}} \quad (21)$$

$$d'_{Tmin} = \min\{d'_T(i, j) | i \in P, j \in Q\} \quad (22)$$

$$d'_{Tmax} = \max\{d'_T(i, j) | i \in P, j \in Q\} \quad (23)$$

$$d'_T(i, j) = \text{abs}(T'(i, j) - \bar{T}'(i, j)) \quad (24)$$

313 Where $d'_T(i, j)$ is the absolute value of difference between the points $T'(i, j)$ and its local
 314 mean value $\bar{T}'(i, j)$ computed over a local overlapping window of size $w \times w$ around the
 315 $T'(i, j)$ after the β – enhancement operation. Now, the membership values in the set T
 316 become more distinct and have high contrast.

318 2.7. Adaptive α and β selection

319 The α and β are the two important parameters for indeterminacy reduction in a neu-
 320 trosophic set. These two parameters can affect the segmentation results. So to determine
 321 the parameters adaptively depending on the characteristics of individual image, we follow

322 the same strategy as [36]. The parameters α and β are computed based on entropy EnI of
 323 subset I as:

$$EnI = - \sum_{i=1}^P \sum_{j=1}^Q pb(i, j) \log_2 pb(i, j) \quad (25)$$

324 where $pb(i, j)$ is the probability of an pixel value at (i, j) in the subset I .

$$En_{max} = \log_2 PQ \quad (26)$$

$$\alpha = \alpha_{min} + \frac{(\alpha_{max} - \alpha_{min})(EnI - En_{min})}{(En_{max} - En_{min})} \quad (27)$$

$$\beta = 1 - \alpha \quad (28)$$

327 Where $P \times Q$ is the dimension of the image. En_{max} and En_{min} are the maximum and
 328 minimum entropy value of I . The α varies in the range $[\alpha_{min} \alpha_{max}]$. The true membership
 329 values of the neutrosophic image after $\alpha - mean$ and $\beta - enhancement$ operations become
 330 the features for segmentation.

331 2.8. Unsupervised feature selection:

332 The unsupervised feature selection is used to select the compact set of significant features
 333 which has minimum correlation between them. Partitioning of the features is done using
 334 the feature similarity Maximal information compression index [38].

335 The Maximal Information compression index (λ_2) for two features x and y is given by

$$2\lambda_2(x, y) = \frac{var(x) + var(y) - \sqrt{(var(x) + var(y))^2 - 4var(x)var(y)(1 - \rho(x, y))^2}}{2} \quad (29)$$

336 Where $var()$ and $\rho()$ denote the variance of a variable and correlation coefficient between
 337 two variables respectively. The (λ_2) is the eigenvalue for the direction normal to the principal
 338 component direction of feature pair (x, y) . The features are partitioned based on k-NN
 339 principle using the feature similarity. The compact set of features are chosen based on λ_2
 340 which are the representative of k-neighbouring features. The value of (λ_2) is zero when the
 341 features are linearly dependent and increases as the amount of dependency decreases.

342 2.9. $\gamma - k -$ means clustering algorithm for NS domain:

343 We use $\gamma - k -$ means clustering [36] for segmentation in NS domain. The clustering algo-
 344 rithm is applied on the true subset of NS image after the $\alpha - mean$ and the $\beta - enhancement$
 345 operations. Let the true subset and the indeterminate subset of NS image after these two
 346 operations become $T_{\alpha\beta}$ and $I_{\alpha\beta}$ respectively. Considering the effect of indeterminacy, the
 347 true subset $T_{\alpha\beta}$ is transformed into a new subset X for clustering as follows.

$$X(i, j) = \begin{cases} T_{\alpha\beta}(i, j) & \text{if } I_{\alpha\beta}(i, j) \leq \gamma \\ \bar{T}_{\gamma}(i, j) & \text{if } I_{\alpha\beta}(i, j) > \gamma \end{cases} \quad (30)$$

348 where $\bar{T}_{\gamma}(i, j)$ is calculated over a local overlapping window of size $w \times w$ around each
 349 $T_{\alpha\beta}(i, j)$. The objective function for the clustering is defined by

$$J_{TC} = \sum_{l=1}^k \sum_{i=1}^H \sum_{j=1}^W \|X(i, j) - Z_l\|^2 \quad (31)$$

350 where k is the number of clusters and $H \times W$ represents the dimension of X . Since, in the
 351 proposed method the number of segments is two, one for text region and other non-text
 352 region, here $k = 2$.

$$Z_l = \frac{1}{n_l} \sum_{X(i, j) \in C_l} X(i, j) \quad (32)$$

353 where J_{TC} is a compactness measure, n is the number of data to be clustered, and C_l is the
 354 l th cluster.

355 3. Proposed Methodology

356 It is well known that in a mixed text document, the text regions and non-text/graphics
 357 regions have distinctly different texture characteristics. With this idea in mind, we extract
 358 the texture features by decomposing it into shearlet sub-bands using DST. The local energy
 359 of each shearlet coefficient is computed using a small overlapping window of adaptive size
 360 around the each coefficient [39]. In order to reduce the uncertainty present in the sub-
 361 bands, each of them is then mapped into NS domain. This is followed by feature selection
 362 and feature dimensionality reduction. The final feature set thus computed is used for feature

363 vector generation and segmentation. The algorithm is illustrated in the Figure 3 and the
 364 steps are explained below. In the figure

365

- 366 • B_1, B_2, \dots, B_n are the shearlet sub-band in DST.
- 367 • T_i, I_i, F_i are the true, indeterminate and false value of B_i where $i = 1, 2, \dots, n$.
- 368 • T'_i is the updated true value [Section 2.6] of shearlet B_i after the iterative α – mean and
 369 β – enhancement operations.
- 370 • Fe_1, Fe_2, \dots, Fe_m are the features after unsupervised feature selection.

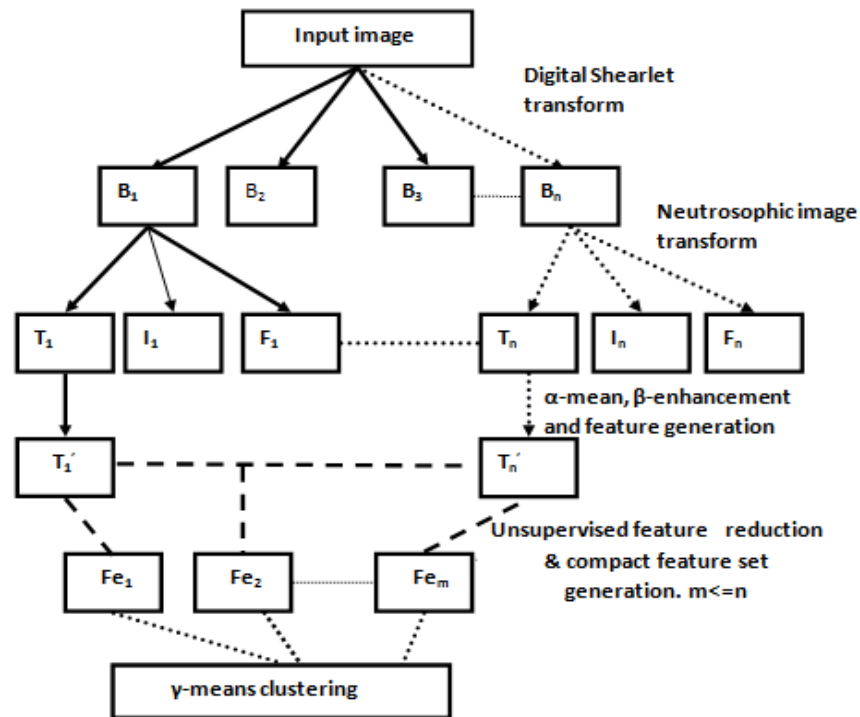


Figure 3: The schematic diagram of the proposed method.

370 **Step 1 : DST transform of the image** - In this step, the input image is transformed into
 372 DST shearlet sub-bands. Each sub-band contains the different scale and directional
 373 information of the image.

375 **Step 2 : Local energy estimation and smoothing of the sub-bands:-** In each sub-
 376 band, the local energy of a shearlet coefficient is computed by a nonlinear operation,
 377 calculated over a small adaptive overlapping window of size w around the coefficient.
 378 Then, the very low energy contained in a shearlet sub-band is removed by a Gaussian
 low-pass (smoothing) filters.

379 **Step 3 : Mapping of the shearlet sub-bands to neutrosophic image:-** To handle the
 380 uncertainties in the sub-bands, each of them is mapped into the neutrosophic image.

381 **Step 4 :The α – mean and the β – enhancement operation on neutrosophic images:-**
 382 In this step, to reduce the uncertainties in the neutrosophic images, each of them
 383 is subjected to the α -mean followed by the β -enhancement operation. These two
 384 operations are done repeatedly until the indeterminate subset entropy of the NS image
 385 in Eq 25 becomes unchanged. That means these two operations are done up to i th
 386 iteration until $|EnI(i) - EnI(i + 1)| \leq \xi$, where ξ is a small positive quantity. The
 387 values of α and β are chosen adaptively(as described in sec 2.7). The true membership
 388 value $T(i,j)$ of each neutrosophic image represents the textures at different scales and
 389 forms the feature for segmentation. The uncertainties within a sub-band is minimized
 390 in this step.

391 **Step 5 : Unsupervised feature reduction and feature vector generation :-** As the
 392 number of features is comparatively large, compaction of the feature set is necessary
 393 in order to achieve higher recognition accuracy and better computation efficiency. To
 394 get a compact set of features, the features are reduced and selected on the basis of
 395 unsupervised feature similarity(described in the section 2.8). The selected features
 396 thus used for generating the multidimensional feature vector. The uncertainties due
 397 to redundant information generation is reduced in this step.

398 **Step 6 :The γ – k – means clustering for segmentation:-** Finally the γ – k – means
 399 clustering (described in section 2.9) is used for segmenting the feature vectors in two
 400 classes, comprising of text and non-text regions. The clustering reduces the uncertain-
 401 ties during segmentation.

402 3.1. Computational complexity of the proposed method :

403 The analysis of computational complexity (worst case) involves the following steps (1)
 404 Computation of R number of $P \times Q$ DST shearlet sub-bands upto scale j where $R =$
 405 $\sum_{j=0}^{nscales-1} 2^{\lceil j/2 \rceil + 2}$ (2) NS image generation and $\alpha - mean$, $\beta - enhancement$ operation (3)
 406 Unsupervised feature selection (4) $\gamma - K - means$ clustering. The overall complexity of the
 407 proposed method is $O(R \times P^2Q^2)$, where the image size is $P \times Q$.

408 4. Experimental results and discussion.

409 4.1. Experimental setup and performance measure :

410 The proposed text region segmentation method combines the merits of multi-resolution
 411 analysis of DST and uncertainty representation using the NS. Experiments were done rigor-
 412 ously, and extensively to judge the ability of the proposed method. The importance of the
 413 parameters α and β for text region segmentation was also studied. It is to be noted that all
 414 the experiments were carried out without a prior knowledge about the input images. The
 415 performance of the proposed method was extensively compared qualitatively and quanti-
 416 tatively with some existing text region segmentation methods. The proposed method was
 417 performed using MATLAB2013 with Pentium IV processor. For the proposed method, the
 values of different parameters are given in the table 1.

Table 1: Values of different parameters used in the proposed method

<i>Parameter</i>	<i>Value</i>
w	5
En_{min}	0
α_{min}	0.01
α_{max}	0.1
ξ	0.001
γ	0.5

418

419 To judge the quantitative performance, four standard measures such as Recall R , Precision
 420 P and F-measure f were used [19]. The performance evaluation was based on Intersection-
 421 over-Union (IOU), with a threshold of 50%, following the standard practice in object recog-
 422 nition [40]. The measures are defined in Eq 33, Eq 34, and Eq 35 respectively [19].

$$R = \frac{|n1|}{|n2|} \quad (33)$$

$$P = \frac{|n1|}{|n3|} \quad (34)$$

$$f = \frac{2 \times P \times R}{(P + R)} \quad (35)$$

425 Where $n1$ is the set of true positive detections, $n2$ is set of the ground truth rectangles and
 426 $n3$ is the set of estimated rectangles.

427 4.2. Datasets:

428 To verify the efficiency of the proposed method, in the present experiment, we used the
 429 document images from the ICDAR2015 [41], ICDAR 2003 [42], ICDAR2011, ICDAR2013,
 430 KAIST [43] and also the scanned images from newspaper, magazine, advertisement found on
 431 the publicly available websites. The size of the images varied from 120×120 to 700×700 .

432 For quantitative performance measure, different methods used different datasets and
 433 some of them are not publicly available. Also, the protocols for measuring the performances
 434 were different. So we have to apply them on the same dataset and use the same protocol for
 435 performance measure to make the comparisons fair. For quantitative performance measure,
 436 we applied all of them to the ICDAR2015 born digital dataset which contains 141 images
 437 and used the protocol in [40] for performance measure. For the text region segmentation
 438 task the ground truth data is provided in the dataset in terms of word bounding boxes.
 439 The reason for applying on this dataset is that born-digital images are inherently low reso-
 440 lution [44]. So ambiguities in the images are more and they are difficult for segmentation.
 441 This is the motivation for using this dataset and we used the updated version of it. For
 442 measuring the performance under different perturbations we used the same dataset. We
 443 also applied our method on KAIST dataset for the quantitative performance measure. It

Table 2: Average performance comparison of segmentation results obtained by different methods using gray level test images from ICDAR2015(Born digital)

<i>Methods</i>	<i>R</i>	<i>P</i>	<i>f</i>
<i>Acharyya</i> [25]	0.58	0.60	0.59
<i>Kumar</i> [26]	0.61	0.63	0.62
<i>Gomez</i> [10]	0.72	0.60	0.66
<i>Maji</i> [28]	0.80	0.82	0.81
<i>Proposed</i>	0.83	0.85	0.84

444 contains 395 test images for text region segmentation.

445

446 4.3. Experiments on gray document images :

447 To show the effectiveness of the proposed method on gray document images, apart from
 448 the standard dataset described earlier, the method was applied to different types of gray
 449 images which include images with non-overlapping text, overlapping text and graphics, the
 450 text of different sizes and orientations.

451 We applied the DST on gray level image. The image was decomposed into 4 level with
 452 one low and three high-frequency shearlets by DST. The three high frequency shearlets were
 453 decomposed into $8(= 2^{\lceil 1/2 \rceil + 2})$, $8(= 2^{\lceil 2/2 \rceil + 2})$ and $16(= 2^{\lceil 3/2 \rceil + 2})$ shearlet sub-bands by the
 454 DST. So total $33(= 1 + 8 + 8 + 16)$ sub-bands of the same size of the input image were gen-
 455 erated. The high-frequency shearlets contained the local information of text and non-text
 456 regions in 8,8 and 16 orientations i.e at each 45° , 45° and 22.5° anti-clockwise directions for
 457 the levels 1,2 and 3. This was followed by local energy estimation of DST coefficients. To
 458 reduce the uncertainty present in the sub-bands and due to redundant information genera-
 459 tion, they were mapped into NS domain for segmentation as described in Section 3.

460 We compared the performance of the proposed method with four published methods.
 461 They are Acharyya [25], Kumar [26], Maji [28] and Gomez [10]. The qualitative results of the
 462 text region segmentation by five different methods (Acharyya [25], Kumar [26], Maji [28],



Figure 4: Column wise (a) Original test images (b) Ground truths for segmentation result (c) Results using Acharyya [25] (d) Results using Kumar [26] (e) Results using Gomez [10] (f) Results using Maji [28] (g) Results using proposed method

463 Gomez [10] and the proposed) are shown in the Figure 4. We applied the proposed method,
 464 and all the other methods considered here on the gray version of the ICDAR2015 dataset of
 465 Born Digital images. The average quantitative performances of the five methods are shown
 466 in the table 2.

467 It is observed from both the qualitative and quantitative measures that on an average the
 468 performance and accuracy of the proposed method are much better than that of Acharyya
 469 [25], Kumar [26], Gomez [10] and Maji [28]. The DST used in the proposed method leads to
 470 a unified treatment of the continuum as well as the digital realm, while still providing opti-
 471 mally sparse approximations of anisotropic features. Thus, the generated features from the
 472 multi-resolution and multi-directional DST shearlet sub-bands and neutrosophic approach
 473 were more effective than the conventional MSERs [10] based method. Moreover, the features
 474 were quite useful in finding the accurate segmentation boundary than the methods in [25]
 475 and [26] which used M-band wavelet and matched wavelet respectively. The M-band wavelet
 476 and matched wavelet have less capacity to capture the edges and curvatures of text regions
 477 than the DST sub-bands. Again, the method in [28] used the multi-resolution analysis
 478 with rough fuzzy clustering for handling the uncertainty during segmentation. However, we
 479 reduced the uncertainties in DST feature in two stages. In the first stage, the uncertainties
 480 in the shearlets were reduced using iterative $\alpha - mean$ and $\beta - enhancement$ operation
 481 during feature generation. That means, the contrasts of the coefficients in the shearlets
 482 were increased and that made them suitable for segmentation. This was followed by feature
 483 set reduction to handle the uncertainties due to redundancy. In the second stage, the un-
 484 certainties during clustering were handled by $\gamma - k - means$ in the NS domain. These two
 485 steps reduction of the uncertainties were more effective than uncertainty reduction in [28].

486 The unsupervised feature selection algorithm in the proposed method takes into account
 487 the uncertainties due to feature redundancy and yields a more compact subset of features.
 488 We tested the effect of unsupervised feature selection in NS domain and the effect of param-
 489 eter k (described in the section 2.8) on the performance. The analytical result is shown in
 490 Appendix A.1.

491 4.4. Experiments on noisy, rotated and dynamic gray level changed document images:

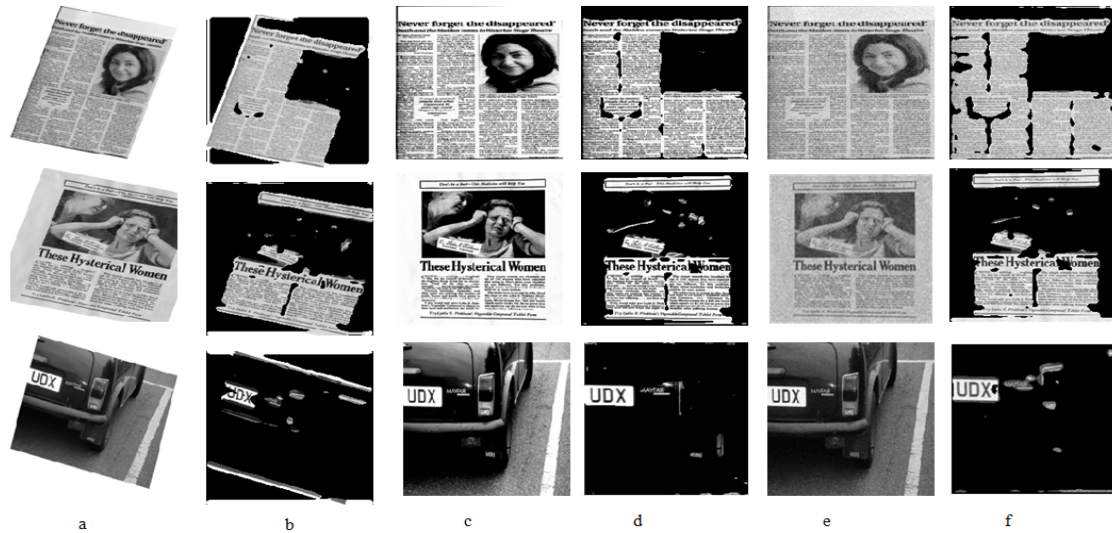


Figure 5: Column wise (a) Rotation of test images by 15.5 degrees clockwise (b) Segmentation results by the proposed method (c) Change by +55 for gray level dynamic range of the test images (d) Segmentation results by the proposed method (e) Test images with Gaussian noise of mean 0 and std 0.16 (f) Segmentation results by the proposed method

492 To demonstrate the robustness of the proposed method, it was applied to a new set of
 493 images with rotation, gray level dynamic range shift, and noise corruption. Some of the test
 494 images and the corresponding text region segmentation results are shown in the Figure 5.
 495 The test images were generated by moderate rotations, adding Gaussian noise and dynamic
 496 gray level intensity change. We compared the quantitative performance of the proposed
 497 method and other text region segmentation methods considered here over noisy, rotated
 498 and dynamic gray level changed images. ICDAR2015 dataset was used for this purpose.

499 The uncertainties in the DST shearlet features increase when noise is added to an image.
 500 So the noise addition affects the text-region segmentation result. The average F-measures'
 501 against average signal to noise ratio (SNR) for all the methods are reported in the graph at
 502 Figure 6. The figure illustrates that our method is robust against a moderate amount of
 503 noise corruption than that of the other methods compared here. The proposed method is ro-
 504 bust against a moderate noise corruption due to the noise handling capacity of shearlets and

505 due to the adaptive α – *means* and β – *enhancement* operations in the proposed method.
 506 These two operations make the noisy pixels in the neutrosophic image homogeneous with
 507 the neighboring pixels and the noise is reduced. Therefore, the uncertainties in the features
 508 get minimized.

509 In the proposed method the DST produces the rotation invariant features due to its ex-
 510 cellent directional property. Additionally, the neutrosophic set theoretic approach reduces
 511 the spatial ambiguity present in the features due to the rotation. As a result, the method
 512 becomes robust against rotation. The average F-measures against the angles of rotation is
 513 shown in the graph at Figure 7. From the graph, it can be said that the proposed method
 514 can also detect non-horizontal text regions more efficiently than the methods compared here.

515 The dynamic gray level modification changes the image pixel intensity, and as a result,
 516 the gray level ambiguity is increased. The ambiguity is successfully reduced in the NS do-
 517 main. The reason is that there will be negligible change in $T(i, j)$ membership values due
 518 to the changes in r values in the Eq. 7. The average F-measures against the positive and
 519 negative gray level shifts are shown in the graphs at Figure 8 and Figure 9 respectively. It
 520 is observed from the graphs that the proposed method is fairly tolerant of moderate changes
 521 in gray level dynamic range.

522

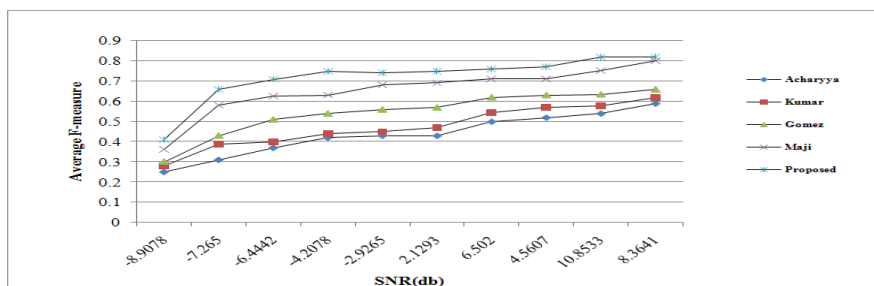


Figure 6: Average F-measure at different SNR on ICDAR2015 dataset(Born digital)

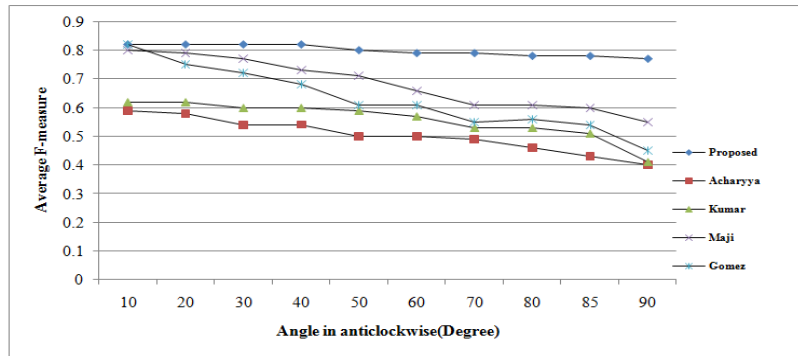


Figure 7: Average F-measure at different orientations on ICDAR2015 dataset (Born digital)

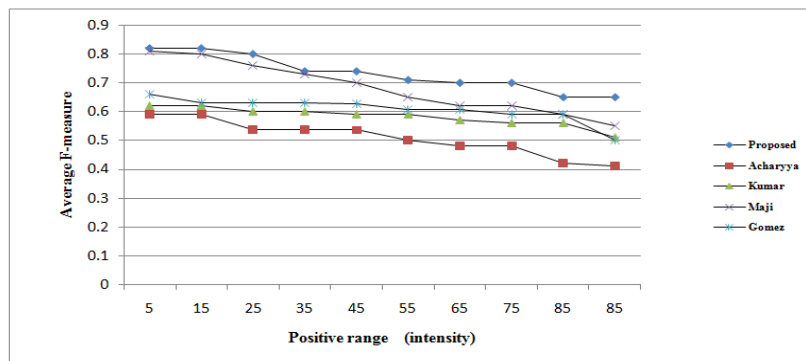


Figure 8: Average F-measure at different positive dynamic range change on ICDAR2015 dataset (Born digital)

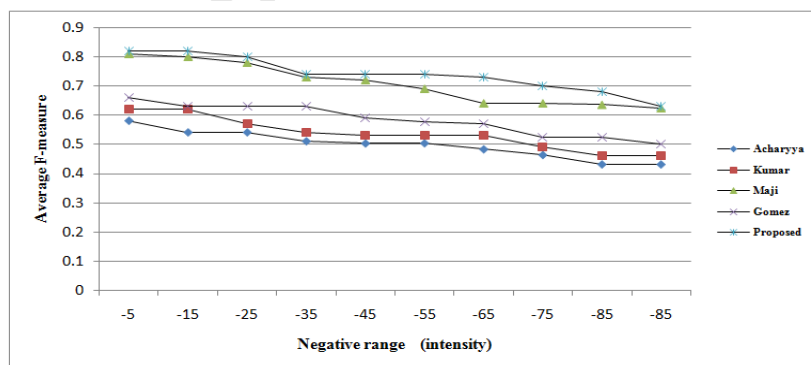


Figure 9: Average F-measure at different negative dynamic range change on ICDAR2015 dataset (Born digital)

523 *4.5. Experiments on colored document images :*

524 To judge the performance of the proposed method on colored text documents, we ap-
 525 plied it to the colored image of the scanned newspaper, magazine etc. For this, at first,
 526 the colored document image was transformed into the YCbCr plane. The YCbCr basically
 527 decouples the intensity and color information, and this representation is very close to the
 528 human perceptual model. The human visual system is less sensitive to chrominance than
 529 luminance [45]. Accordingly, the weights of features generated from the above three planes
 530 were based on the convention perceptual importance, as used in the JPEG 2000 that is Y:
 531 Cb: Cr=4:2:1. Each of the Y, Cb, Cr color plane was transformed into a set of frequency
 532 shearlets using DST. The total number of shearlet sub-band generated was $99(= 3 \times 33)$.
 533 The sub-bands were then mapped to NS domain for uncertainty reduction. This was fol-
 534 lowed by feature vector generation and segmentation as described in section 3.

535 We compared the performance of the proposed method with three published methods in
 536 [29], [36] and [10] on color documents. The method in [36] used dyadic wavelet transform
 537 and NS set theocratic approach for color texture segmentation. As the basic assumption
 538 of the proposed method is that the text and non-text regions are two different textures,
 539 the proposed method is also capable of segmenting any color texture image which contains
 540 two different textures. Here it was assumed that the method in [36] was also capable of
 541 differentiating the text and non-text regions.

542 The quantitative performance of all the three methods and the proposed method were
 543 computed using three different metrics described in Eq 33, Eq 34 and Eq 35 respectively.
 544 The color text region segmentation results due to four different methods ([29], [36], [10]
 545 and the proposed) are shown in the Figure 10. For quantitative performance, we applied all
 546 the four methods and the method by Cho [16] on the color images of ICDAR2015. The
 547 results are shown in table 3. We also compared the quantitative performance of the pro-
 548 posed method with the methods in [10] and [15] on KAIST dataset. The performances are
 549 shown in table 4. It is observed both from the qualitative and quantitative results that the
 550 accuracy of the proposed method on average is much better than that of the other methods.

551 In the proposed method for a color document image, both the color and texture infor-

552 mation were combined. The DST shearlet features are able to extract the texture features
 553 more accurately than dyadic wavelet transform. Hence, the proposed method provided the
 554 better result than the method in [36]. Again, from the performances, it can be said that the
 555 feature set and the uncertainty handling technique of the proposed method during feature
 556 generation is better than that of [29] during segmentation. The method in [29] used fuzzy
 557 c-means for segmentation. The proposed method also performed better than that of the
 558 seed based method in [15]. The reason may be that the method in [15] depended on the
 559 Canny edge detector for generating the seed. The detector may not work well in the complex
 560 background and it can not handle the uncertainties. Again the method in [16] combined
 561 the MSERs and Canny edge detector to detect the text regions, which are less efficient than
 562 the proposed anisotropic features of the text captured by the DST.

Table 3: Average performance comparison of text region segmentation results obtained by different methods using color test images from ICDAR2015(Born digital)

<i>Methods</i>	<i>R</i>	<i>P</i>	<i>f</i>
<i>Kundu</i> [29]	0.67	0.75	0.71
<i>Sengur</i> [36]	0.74	0.79	0.77
<i>Gomez</i> [10]	0.73	0.76	0.74
<i>Cho</i> [16]	0.77	0.83	0.79
<i>Proposed</i>	0.83	0.84	0.83

Table 4: Average performance comparison of text region segmentation results obtained by different methods using color test images from KAIST

<i>Methods</i>	<i>R</i>	<i>P</i>	<i>f</i>
<i>Gomez</i> [10]	0.78	0.66	0.71
<i>Bai</i> [15]	0.89	0.83	0.86
<i>Proposed</i>	0.89	0.86	0.87



Figure 10: Column wise (a) Original color test image (b) Ground truth for test color image segmentation (c) Results using Kundu. [29] (d) Results using Sengur [36] (e) Results using Gomez [10] (f) Result using proposed method

563 4.6. Comparison with deep learning based methods

564 For the learning based method, we used the CNN with the shearlet features in NS domain
565 for text region detection. We trained the network by using the reduced features in the NS
566 domain. For classification using CNN, we followed the same protocol as used in [46] for
567 text region segmentation using multi-resolution analysis. The CNN text detector had 5
568 layers. They were 2 convolution layers, 2 average pooling layers, and a fully connected
569 layer. The output layer was consisted of 2 nodes which were text and non-text. To train the
570 CNN we used the 229 training image from ICDAR2011 and 258 images from ICDAR2005
571 dataset. We synthetically created the text and non-text patches from them. We compared
572 our method with state-of-the-art deep learning based text detection methods. We tested
573 on the same dataset of ICDAR2011 and ICDAR2013 for comparison which contains 255
574 and 233 test images respectively. In Huang’s [18] and He’s [17] deep learning method, the
575 MSERs features were used as input to the CNN. The quantitative results are shown in Table
576 5. From the result, it can be said that the proposed method performed better than that
577 of He and Huang. The reason is that the MSERs features used in the deep learning, are
578 less robust in the rotation than DST shearlet features used in the proposed method. Hence
579 the MSER features become more uncertain than that of the DST features. Moreover, the
580 uncertainties in the DST features were reduced in NS domain before entering into CNN,
581 where no such provisions were available in their methods.

Table 5: Experimental results of Deep learning based methods using ICDAR dataset

<i>Methods</i>	<i>Dataset</i>	<i>R</i>	<i>P</i>	<i>f</i>
<i>He</i> [10]	<i>ICDAR2011</i>	0.74	0.91	0.82
<i>Huang</i> [15]		0.71	0.88	0.78
<i>Proposed</i>		0.74	0.94	0.84
<i>He</i> [10]	<i>ICDAR2013</i>	0.73	0.93	0.82
<i>Proposed</i>		0.75	0.94	0.83

582

583 4.7. Additional qualitative results

584 Some additional results by the proposed method from ICDAR2015(Born digital) dataset
 585 and KAIST dataset are shown in Figure11. Further, the method was applied to some more
 586 document image in different orientations. In these images, ambiguities occurred due to
 587 different orientations. Some test image data from Computer vision Laboratory [47] and
 588 corresponding segmentation results are shown in Figure 12. The visual inspection again
 589 indicates that the method is fairly tolerant of different orientations .

590



Figure 11: Row wise (a) Row1 shows images from ICDAR2015 and KAIST. Row2 shows text region segmentation results of corresponding images in the same sequence, using the proposed method.

591 4.8. Role of α and β parameters :

592 In this sub-section, we studied the role of the two parameters α and β for text region
 593 segmentation. As already stated, the operations are used to reduce indeterminacy in the NS
 594 domain. If these two parameters are not appropriate, the proposed method may not generate
 595 features properly for text and non-text regions. This may lead to poor segmentation result.
 596 In Figure 13, three results are shown for different values of α and β . When the values of
 597 α and β were both 0.83 the result is shown in the Figure 13(b). The result in Figure 13(c)
 598 was obtained when values of α , β were 0.75 and 0.25 respectively. But when the values
 599 were computed adaptively based on the statistics of the I subset, we got the most consistent
 600 result as shown in Figure 13(d). The reason is that α is high(low) and β is low(high) when



Figure 12: Column wise (a) Color text image set from Computer vision Laboratory dataset (b) Text region segmentation results of corresponding images in (a) in the same sequence, using the proposed method. (c) Another set of Color text image set corresponding to the same row in (a) with the different orientation (d) Text region segmentation results of corresponding images in (c) in the same sequence, using the proposed method.

601 entropy of I is high(low). High entropy means the image contains more indeterminacy i.e
 602 it is difficult to decide whether a pixel is text or non-text pixels. So shearlet sub- band
 603 with high indeterminacy i.e $I \geq \alpha$ should be made homogenous by α -means with contrast
 604 enhancement by β -enhancement to generate the proper feature vector. In the first two cases,
 605 the values of α and β were chosen arbitrarily without considering the document image in
 606 hand and hence could not generate the proper feature vector.

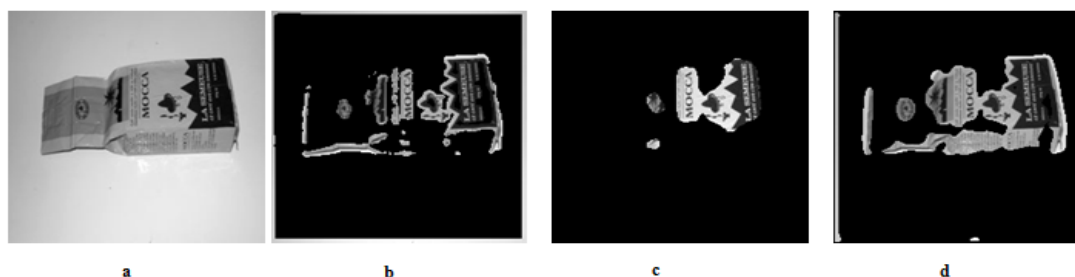


Figure 13: Column wise (a) Original document test image (b) Results using $\alpha = 0.83$ and $\beta = 0.83$
 (c) Results using $\alpha = 0.75$ and $\beta = 0.25$ (d) Results using adaptive α and β

607

608 5. Conclusion.

609 In this paper, we have proposed an DST and NS based text region segmentation method-
610 ology in the complex document images. The method judiciously combines the advantages
611 of multi-scale and multi-directional shearlet feature extraction tool like DST together with
612 the uncertainty handling capability of the neutrosophic set. This is done in order to achieve
613 the segmentation result with higher accuracy in comparison to many existing techniques
614 reported in the literature. It is also to be noted that the proposed method performs equally
615 well in another type of segmentation like two class color and gray level texture segmentation.
616 This method also shows robustness in performance under different perturbations like mod-
617 erate rotation, dynamic range changes, and noise corruption. With suitable modification in
618 the current method, it can be extended for the multiclass problem, which is being currently
619 investigated.

620 6. Acknowledgement

621 We would like to thank Dr. Weilin Huang for providing us the code of their paper [17].

622 7. Reference

- 623 [1] C. A. Murthy, S. K. Pal, Histogram thresholding by minimizing gray level fuzziness, *Information*
624 *Sciences* 60 (1992.) 107–135.
- 625 [2] A. Rosenfield, Fuzzy geometry: An updated overview, *Information Sciences* 110 (1998.) 127–133.
- 626 [3] S. K. Pal, A. Ghosh, M. K. Kundu, *Soft Computing for Image Processing*, Springer-Verlag Berlin
627 Heidelberg GmbH, 1st edition, 2000.
- 628 [4] G. Nagg, S. Seth, M. Viswanathan, A prototype document image analysis system for technical
629 journals, *Computer*. 25 (1992.) 10–22.
- 630 [5] D. Chen, J.-M. Odobez, H. Bourlard, Text detection and recognition in images and video frames,
631 *Pattern Recognition* 37 (2004) 595–608.
- 632 [6] S. M. Lucas., ICDAR2005 text locating competition results ., *Proceedings of International Conference*
633 *on Document Analysis and Recognition*. 1 (2005.) 80–84.
- 634 [7] E. Haneda, C. Bouman, Text segmentation for mrc document compression, *IEEE Transactions on*
635 *Image Processing* 20 (2011) 1611–1626.

- 636 [8] C. Yi, Y. Tian, Text string detection from natural scenes by structure-based partition and grouping.,
637 IEEE Transactions on Image processing. 20 (2011.) 2594–2605.
- 638 [9] J. Matas, O. Chum, M. Urban, T. Pajdla, Robust wide baseline stereo from maximally stable external
639 regions., Proceedings of British Machine Vision Conference (2002) 384–396.
- 640 [10] L. Gomez, D. Karatzas, Multi-script text extraction from natural scenes., Proceedings of International
641 Conference on Document Analysis and Recognition. (2013.) 467–471.
- 642 [11] X. Yin, X. Yin, H. H. K. Hung, Robust text detection in natural scene images, IEEE Transactions on
643 Pattern Analysis and Machine Intelligence. 36 (2014.) 970–983.
- 644 [12] A. Gonzalez, L. Bergasa, J. Yebes, S. Bronte, Text location in complex images., Proceedings of
645 International Conference on Pattern Recognition. (2012.) 617–620.
- 646 [13] C. Shi, C. Wang, B. Xiao, Y. Zhang, Scene text detection using graph model built upon maximally
647 stable extremal regions., Pattern Recognition letter . 34 (2013.) 107–116.
- 648 [14] Y. Li, H. Lu, Scene text detection via stroke width, Proceedings of International Conference on Pattern
649 Recognition (2012) 681–684.
- 650 [15] B. Bai, F. Yin, C. L. Liu, A seed-based segmentation method for scene text extraction., IAPR
651 International Workshop on Document Analysis Systems (2014) 262–266.
- 652 [16] H. Cho, M. Sung, B. Jun, Canny text detector: Fast and robust scene text localization algorithm.,
653 International Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 3566–3573.
- 654 [17] T. He, W. Huang, Y. Qiao, J. Yao, Text-attentional convolutional neural network for scene text
655 detection., IEEE Transactions on Image Processing 25 (2016) 2529–2541.
- 656 [18] W. Huang, Y. Qiao, X. Tang, Robust scene text detection with convolution neural network induced
657 MSER trees, Proceedings of European Conference on Computer Vision LNCS 8692 (2014) 497–511.
- 658 [19] Y. Zhu, C. Yao, X. Bai, Scene text detection and recognition: recent advances and future trends,
659 Frontiers of Computer Science 10 (2016) 19–36.
- 660 [20] G. Zhou, Y. Liu, Q. Meng, Detecting multilingual text in natural scene., Proceedings of International
661 Symposium on Access spaces . (2011.) 116–120.
- 662 [21] M. Zhao, S. Li, J. Kwork, Text detection in images using sparse representation with discriminative
663 dictionaries., Image and Vision Computing . 28 (2010.) 1590–1599.
- 664 [22] C. Liang, P. Y. Chen, Dwt based text localization., International Journal of Applied Science Engineer-
665 ing. 2 (2004.) 105–116.
- 666 [23] W. Chan, G. Coghill, Text analysis using local energy., Pattern Recognition . 34 (2001.) 2523–2532.
- 667 [24] S. Nirmal, P. Nagabhushan, Foreground text segmentation in complex color document images., Signal
668 Image and Video Processing . 6 (2012.) 669–678.
- 669 [25] M. Acharyya, M. Kundu, Document image segmentation using wavelet scale-space features., IEEE

- 670 Transactions on circuits and systems on video technology. 12 (2002.) 1117–1127.
- 671 [26] S. Kumar, R. Gupta, N. Khanna, S. Chaudhury, S. Joshi, Text extraction and document image
672 segmentation using matched wavelets and mrf model., IEEE Transactions on Image processing . 16
673 (2007.) 2117–2128.
- 674 [27] S. Roy, M. Kundu, G. .Granlund, Uncertainty relations and time-frequency distributions for unsharp
675 observables, Information Sciences 89 (1996.) 193–209.
- 676 [28] P. Maji, S. Roy, Rough -fuzzy clustering and multiresolution image analysis for text-graphics segmen-
677 tation., Applied soft computing. 30 (2015.) 705–721.
- 678 [29] M. Kundu, S. Dhar, M. Banerjee, A new approach for segmentation of image and text in natural and
679 commercial text documents., Proceedings of International Conference on Communications ,Devices
680 and Intelligent system. (2012.) 86–88.
- 681 [30] G. Kutyniok, W. Q. Lim, R. Reisenhofer, Shearlab 3d: Faithful digital shearlet transforms based on
682 compactly supported shearlets., Numerical Analysis (math.NA) arXiv:1402.5670 (2014).
- 683 [31] G. Kutyniok, W. Q. Lim, G. Steidl, Shearlets: Theory and applications., GAMM-Mitt 37 (2014)
684 259–280.
- 685 [32] F. Smarandache, A Unifying Field in Logic,Neutrosophy, Neutrosophic Set,Neutrosophic Probability,
686 American Research Press, 3rd edition, 2003.
- 687 [33] Y. Guo, H. Cheng, A new neutrosophic approach to image segmentation, Pattern Recognition . 42
688 (2009.) 587–595.
- 689 [34] H. Cheng, Y. Guo, A new neutrosophic approach to image thresholding, New Mathematics and Natural
690 Computation. 4 (2008.) 291–308.
- 691 [35] Y. Guo, H. Cheng, A new neutrosophic approach to image denoising, New Mathematics and Natural
692 Computation. 5 (2009.) 653–662.
- 693 [36] A. Sengur, Y. Guo, Color texture segmentation based on neutrosophic set and wavelet transform,
694 Computer Vision and Image Understanding. 115 (2011.) 1134–1144.
- 695 [37] G. Kutyniok, W. Q. Lim, X. Zhuang, Digital shearlet transforms. In Shearlets: Multiscale analysis for
696 multivariate data., Springer, 2012.
- 697 [38] P. Mitra, C. Murty, S. Pal, Unsupervised feature selection using feature similarity., IEEE Transactions
698 on Pattern Analysis and Machine Intelligence . 24 (2002.) 301–312.
- 699 [39] M. Kundu, M. Acharyya, M-band wavelets:application to texture segmentation for real life image
700 analysis., International Journal of Wavelets, Multiresolution and Information Processing. 1 (2003.) 115–
701 119.
- 702 [40] M. Everingham, L. V. G. and C. K. I. Williams J. Winn, A. Zisserman, The pascal visual object
703 classes (voc) challenge, International Journal of Computer Vision 88 (2010.) 303–338.

- 704 [41] ICDAR2015 dataset, <http://rrc.cvc.uab.es/> (2015.).
- 705 [42] ICDAR2003, <http://algoval.essex.ac.uk/icdar/Datasets.html> (2003.).
- 706 [43] Kaist scene text database, www.iapr-tc11.org/mediawiki/index.php/KAIST_Scene_Text_Database
707 (2011.).
- 708 [44] D. Karatzas, S. R. Mestre, J. Mas, F. Nourbakhsh, P. P. Roy, ICDAR 2011 Robust Reading Competi-
709 tion., International Conference on Document Analysis and Recognition. (2011) 1485–1490.
- 710 [45] N. Plataniotis, A. Venetsanopoulos, Color image processing and applications., Springer Verlag, Heidel-
711 berg . (2000).
- 712 [46] T. Kobchaisawat, T. H. Chalidabhongse, Thai text localization in natural scene images using convo-
713 lutional neural network., Signal and Information Processing Association Annual Summit and Confer-
714 ence(APSIPA) (2014).
- 715 [47] Computer Vision Laboratory dataset, <http://www.vision.ee.ethz.ch/datasets/> (2003.).

716 Appendix A. NS image representation of gray level values of an image

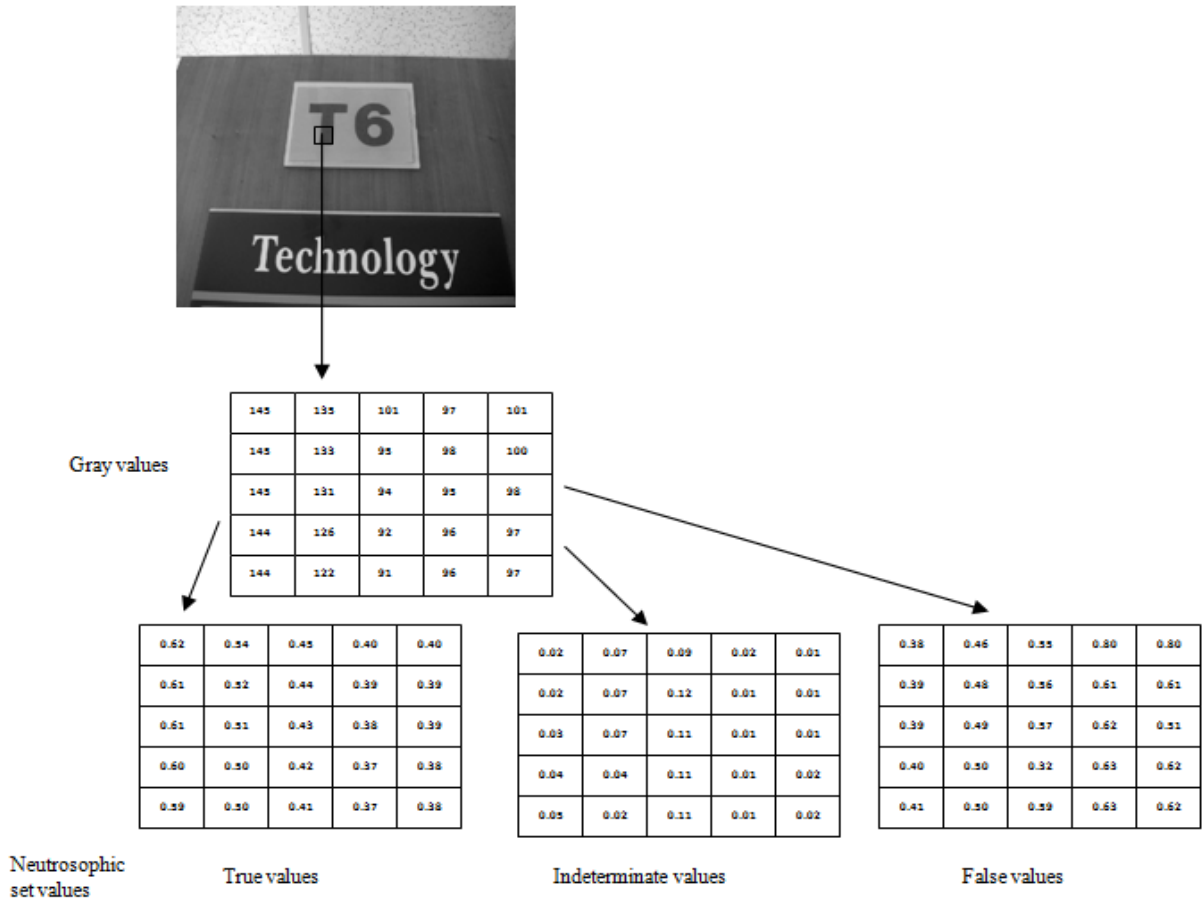


Figure A.14: Neutrosophic set representation of a small portion (marked square) of gray scale image

717 The image in the Figure A.14 shows the gray values of the original image within the
 718 square box and the true, the false and the indeterminate membership values of the corre-
 719 sponding neutrosophic image obtained by the Eq 7-Eq 12. From the third column of the
 720 indeterminate subset matrix, it can be said that pixels in that column are the edge pixels,
 721 and which is correct.

722

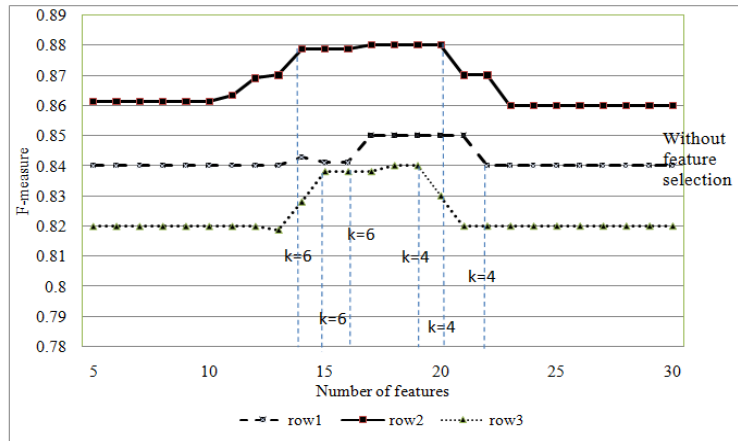


Figure A.15: Variation in f-measure with the size of feature set on three test images. The vertical dotted lines represent the number of features and corresponding f-measure at $k=4$ or $k=6$

723 *Appendix A.1. Importance of unsupervised feature selection :*

724 To examine the importance of unsupervised feature selection, we applied the proposed
 725 method with and without the feature selection, on the gray level images. The quantitative
 726 performance on three images in row 1, row2 and row3 in Figure 4 are illustrated in the
 727 graph shown in Figure A.15. In the figure, the f-measures are shown against the number of
 728 features. It also illustrates the effect of parameter k in the feature selection algorithm. The
 729 figure shows that f-measure becomes low without the feature selection algorithm or when
 730 the number of features is very low. It is observed that the number of features generated is
 731 lower when the value of k is high. In the figure, the vertical dashed lines point the number
 732 of features created when $k=6$ or $k=4$ for each sample image.