

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/312251128>

# Reproducing tables in scanned documents

Article in Journal of the National Science Foundation of Sri Lanka · December 2016

DOI: 10.4038/jnfsr.v44i4.8019

CITATIONS  
3

READS  
98

2 authors:



Akmal Jahan MAC  
South Eastern University of Sri Lanka

21 PUBLICATIONS 85 CITATIONS

[SEE PROFILE](#)



Roshan Ragel  
University of Peradeniya

200 PUBLICATIONS 1,184 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



PhD research of S. M. Vidanagamachchi under the supervision of myself (Dr. S. D. Dewasurendra), Dr. R. G. Ragel and Prof. M. Niranjan [View project](#)



Plagiarism Detection in Text based Documents [View project](#)

## RESEARCH ARTICLE

# Reproducing tables in scanned documents

**M.A.C. Akmal Jahan<sup>1\*</sup> and Roshan G. Ragel<sup>2</sup>**

<sup>1</sup> Postgraduate Institute of Science, University of Peradeniya, Peradeniya.

<sup>2</sup> Department of Computer Engineering, Faculty of Engineering, University of Peradeniya, Peradeniya.

Revised: 20 April 2015; Accepted: 17 March 2016

**Abstract:** Digitisation is a process of representing real world objects in digital format. The rapid conversion of material available in printed form to editable digital form requires a significant amount of work if we are to maintain the format and the style of the electronic documents similar to their printed counterparts. Most of the existing digitisation procedures cover only text, and it has to go beyond OCR (text only) for making the text inside objects such as tables searchable, while preserving the format and converting the objects to an editable form to modify and reprint the content, for which the processes of detection and recognition are important. In past research the table detection process mostly followed certain assumptions such as i) either having rule lines or no lines and ii) Manhattan or multi-column layout, whereas the recognition process assumed that the considered tables have already been detected and fails to preserve their formatting features for future modifications and re-printing. To address these issues, we proposed a simple and fast algorithm using local thresholds for word space and line height, which locates all types of tables and extract their formatting features. From the experiments performed on 353 records, we have achieved a much higher detection ability than the earlier algorithm and is superior as it performs extended layout analysis, elimination of header-footer and bulleted-numbered sections before the reconstruction of tables. While reconstructing the extracted table, the most prominent features of the tables are preserved with the earlier formats. The algorithm has an advantage of linear complexity as it performs the conversion locally.

**Keywords:** Digitisation, format preservation, table detection.

## INTRODUCTION

Complete digitisation of printed materials play a significant role in building digital libraries in which modifying, re-printing, and searching of content are some important tasks. An enormous amount of manual

effort is required to maintain the format and appearance of electronic documents as identical to their printed documents. The focus here is converting printed documents into editable format while preserving both the content and the format of the records.

In the past, most of the optical character recognition (OCR) approaches have focused only on character recognition of the text and extraction of format features of the characters. Nonetheless, most of the printed documents contain not only characters but also associated non-text objects such as tables, charts and graphics. This leads to a challenge in detection and preservation of their features for future editing during the digitisation process. Therefore, the current digitisation has to go beyond OCR (text only) for making the text inside objects such as tables searchable, while preserving the format and making the tables editable to modify and reprint the content. Detection and recognition are the two primary processes involved when handling tables in printed documents. The table detection process mostly follows certain assumptions such as having rule lines or no lines and Manhattan or multi-column layout. Further, the recognition process assumes that the considered tables have already been detected and fails to address the preserving of their formatting features for future modifications and re-printing. Therefore, the objective of this research study was to determine a way to detect, recognise and reconstruct tables from document images to help the process of preserving and reproducing documents with features. As earlier studies have completely addressed the OCR related issues for text (Ajward *et al.*, 2010), we focused on locating different types of tables in documents and treating them with their existing structural features to reconstruct them. Therefore,

\* Corresponding author ([akmaljahan@fas.seu.ac.lk](mailto:akmaljahan@fas.seu.ac.lk))

this study handles tables in three steps: i) locating tables; ii) recognition of tables and iii) reproducing them.

Tables vary in structure from regular text and is composed of ruled lines, decorated lines or no lines. The basic fabric of a table depends not on the line but in its building blocks such as rows, columns and fields. Tables have physical and logical structures (Zanibbi *et al.*, 2003). The physical structure determines the regions of a table and is used in region detection in a document. The logical structure determines the relational information of a table such as integral parts, how they form a table, and helps in table structure recognition. Recently, many researchers have focused on table recognition rather than locating the table, or table detection (Kieninger & Dengel, 2001; Wang *et al.*, 2001; Yildiz *et al.*, 2005; Watanabe *et al.*, 2006; Oro & Ruffolo, 2009; Fang *et al.*, 2011). Most of them assume that the table area is already known, and their work has focused only on the extraction of its logical structure. The existing table recognition systems lack the facility to locate tables in a document and the other non-text objects for reproducing them with the existing formatting. Rarely a few works cover both table detection and recognition (Laurentini & Vaida, 1992; Namboodiri, 2004). The current trend of document analysis recently turned into detection of tables in camera-captured images with different illumination, contrast and distortion (Wonkyo *et al.*, 2015). There is no study in the recent past, which addresses the reproduction of tables with formatting.

Table detection processes can have different approaches based on the input document types such as scanned images and electronic text documents (Yildiz *et al.*, 2005; Oro & Ruffolo, 2009; Fang *et al.*, 2011). The features used for table detection are geometric features such as ruled lines (Hu *et al.*, 2000; Cesarini *et al.*, 2002; Gatos *et al.*, 2005; Kasar *et al.*, 2013), pixel distributions, white gaps (Mandal *et al.*, 2006), header and the trailer pattern of the table (Harit & Bansal, 2012). Most of the past studies have only focused on single column page layout, while a few recent studies have focused on multi-column page layout (Wang *et al.*, 2001; Shafait & Smith, 2010). Namboodiri (2004) addressed the problem of table detection and recognition for online handwritten documents whereas Chen and Lopresti (2011) addressed this issue for off-line handwritten documents. Laurentini and Vaida (1992) used horizontal and vertical ruling lines to identify tables and excluded non-tabular areas from the printed documents. Hu *et al.* (2000) proposed a method, which does not depend on ruling lines and is document independent. Gatos *et al.* (2000) located and reconstructed tables with intersection points from

horizontal and vertical lines. Mandal *et al.* (2006) proposed an algorithm, which assumes the presence of substantially larger gaps between columns and table fields. This system can locate the tables only if they do not contain any ruling lines or are already removed in pre-processing. Also, tables with multi-line headings or heterogeneous placement of cells can result in erroneous detection in this approach. Harit and Bansal (2012) proposed a method for table detection using header and trailer patterns, which depend on rational indications rather than functional meaning. Shafait and Smith (2010), and Smith (2009) have used tab-stop detection for the layout analysis of document images and then used the alignment information of columns for finding tables. Their algorithm does not work when full page tables are present.

From the overall analysis, it was noted that all the previous works handle only a certain page layout, table category with assumptions and consider some definite type of tables rather than all range of tables in different categories. Therefore, we need to identify a novel way to locate all types of tables from the text content of scanned documents with different layouts to reproduce them with their formatting features. The work recently presented by Jahan and Ragel (2014) focused on locating all the different types of tables in printed text. However, their algorithm has some limitations on handling i) header and footer section with rule lines; ii) bullets and numbering points; and iii) multiple column document layout, which finally lowers the performance of the algorithm. Moreover, their algorithm does not represent the way to reproduce a table after locating. For this reason, in this paper we extend the work presented by Jahan and Ragel (2014), and propose an extended algorithm to eliminate the issues. However, the issues identified in previous studies such as heterogeneous placement of cells and full page tables on the page, still need to be resolved. The algorithm suggested in this paper can address the detection of all categories of tables with Manhattan or multi-column layout, with, without or partial ruled lines and a way to reconstruct them. The existing algorithm (Jahan & Ragel, 2014) is improved by applying several modules to eliminate the header-footer section and bullets-numbering parts to purify the algorithm and enhance the throughput.

## METHODOLOGY

The sample document images consist of several parts of the content such as text and non-text objects. The text content consists of regular text lines, headings, a full range of tables with or without rule lines, equations,

**Table 1. Health Status by Selected Characteristics: 2010**  
(Numbers in thousands. Only people in the noninstitutionalized population)

Characteristic	Total number	Health status (percent)					
		Excellent	Very Good	Good	Fair	Poor	Estimate error
All people .....	304,814	32.7	0.8	38.8	0.8	24.1	0.4
Sex .....							
Male .....	149,421	33.9	0.8	35.1	0.3	24.6	0.4
Female .....	155,393	31.6	0.8	32.7	0.3	26.9	0.4
Race and Hispanic origin? .....							
White, non-Hispanic .....	197,798	32.7	0.8	35.4	0.3	24.8	0.4
Black .....	36,154	29.8	0.7	32.1	0.2	25.9	0.6
Other, non-Hispanic .....	21,231	35.3	0.8	31.6	0.7	24.4	0.8
Hispanic .....	49,831	33.8	0.7	32.3	0.6	25.0	1.9
Age .....							
16 to 19 years .....	74,801	32.4	0.6	37.1	0.3	17.7	0.3
20 to 24 years .....	230,012	24.0	0.9	34.8	0.2	28.1	0.9
Less than 65 years .....	265,947	36.2	0.8	35.9	0.3	22.6	0.1
65 years and over .....	86,534	18.1	0.8	34.8	0.4	31.6	0.5
Family income as a percentage of poverty threshold .....							
Less than 100 percent .....	111,541	29.7	0.4	29.8	0.4	26.5	0.2
Less than 200 percent .....	49,508	31.5	0.6	29.8	0.5	25.2	0.3
100 percent or less than 200 percent .....	263,9	30.3	0.8	30.9	0.5	25.6	0.2
200 percent or more .....	190,670	32.9	0.7	32.0	0.4	22.8	0.3
200 percent or less than 300 percent .....	52,197	30.5	0.6	32.5	0.5	26.9	0.4
300 percent or less than 400 percent .....	30,084	32.7	0.6	34.4	0.6	24.9	0.4
400 percent or more .....	96,938	32.4	0.6	36.1	0.4	20.4	0.3

Standard error estimates were calculated using replicate weights. FIPS Method.

1 Federal survey now give respondents the option of reporting more than one race. Therefore, two basic ways of defining a race group are possible. A group, such as Black, may be defined those who reported Black and no other race (the race-as-one or the single-race concept) or as those who reported Black regardless of whether they also reported one or more other races (the race-in-addition concept). The Census Bureau uses the race-in-addition concept. This table shows data for people who reported they were a single race. Use of the single-race concept does not imply that it is the preferred method of reporting or analyzing data. The Census Bureau uses a variety of approaches. In the report, the term "non-Hispanic White" refers to people who are not Hispanic and are white. The term "non-Hispanic Black" refers to people who are not Hispanic and are black. The term "Asian" refers to people who are not Hispanic and reported Asian alone, Pacific Islander alone, American Indian alone, Alaskan Native alone, or multiple races.

The poverty universe is slightly smaller than that reported under "All People" as it excludes people less than 15 years old with no contributing relatives, who are not in the labor force, and institutionalized people.

Source: U.S. Census Bureau, Survey of Income and Program Participation, 2003 Panel, wave 7 (total module and core survey data. For information on confidentiality protection and sampling and nonsampling error, "All People" it excludes people less than 15 years old with no contributing relatives, who are not in the labor force, and institutionalized people.

**Table 1. Source retrieval results with respect to retrieval performance and cost-effectiveness.**

Team	Downloaded sources	Total workload	Workload to 1st detection	No. runtime					
	F <sub>1</sub>	Precision	Recall	Queries Downloads					
Elizalde	0.17	0.12	0.44	44.50	107.22	16.85	15.28	5	241.7 m
Gilliam	0.04	0.02	0.10	16.10	33.02	18.80	21.70	38	15.1 m
Haggag	0.44	0.63	0.38	32.04	5.93	8.92	1.47	9	152.7 m
Kong	0.00	0.00	0.65	48.50	5601.47	2.46	285.66	3	4098.0 m
Lee	0.35	0.50	0.33	44.04	11.16	7.74	1.72	15	310.5 m
Nourian	0.10	0.15	0.10	4.91	13.54	2.16	5.61	27	25.3 m
Suchomel	0.06	0.04	0.23	12.38	261.95	2.44	74.79	10	1637.9 m
Vesely	0.15	0.11	0.35	16.21	81.03	18.00	5.07	16	655.3 m
Williams	0.47	0.55	0.50	116.40	14.05	17.59	2.45	5	1163.0 m

**Figure 1:** Different types of tables in document image: a) table with partially and fully bounding lines; b) tables with parallel lines; c) tables without rule lines

header, footer and page numbers with different font sizes, types, styles and different layouts. The background of the document is white or gray coloured, and there are no graphics or watermarks on the page. The document pages are with Manhattan or multi-column layouts and saved as tiff, JPEG or PNG image format. Around 353 scanned images, the documents with tables, were categorised into three types such as tables with: a) fully and partially bounding lines; b) parallel lines; and c) without rule lines as shown in Figure 1.

## Pre-processing

Pre-processing plays an important role in the detection process. Initially, it involves the process of binarisation in which colour images or grey-scale images are converted to a binary image using adaptive threshold

mechanism (Gatos *et al.*, 2004). The marginal non-textual noisy border is removed, and image enhancement is accomplished by dilation using a structuring element as shown in Figure 2.

## Algorithm

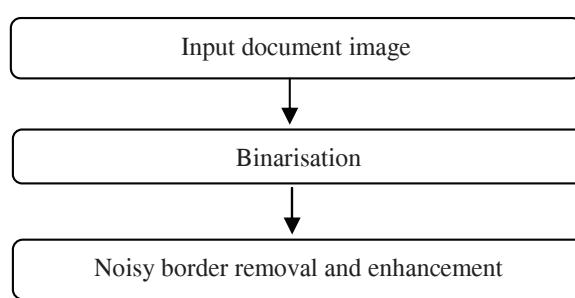
The new algorithm for reproducing the tables is implemented by extending the previous algorithm by Jahan and Ragel (2014) as presented in Figure 3, which follows a number of stages: i) threshold computation for document layout analysis; ii) threshold computation for word space and line height; iii) eliminate header and footer sections; iv) removing bullets and numbering; v) table detection within text content; and vi) recognition and reconstruction of tables with the formatting features.

### Threshold computation for layout analysis

In the initial step, all layout documents have been separated and each column in the page has been treated with the table detection algorithm. Here, we have used a threshold value for inter-column space. The algorithm for layout analysis is given in Figure 4.

### Removing bullets and numbering from the content

The algorithm is improved by eliminating bullets and numbering from the content as they interrupt during threshold calculation of the standard word space. Here,

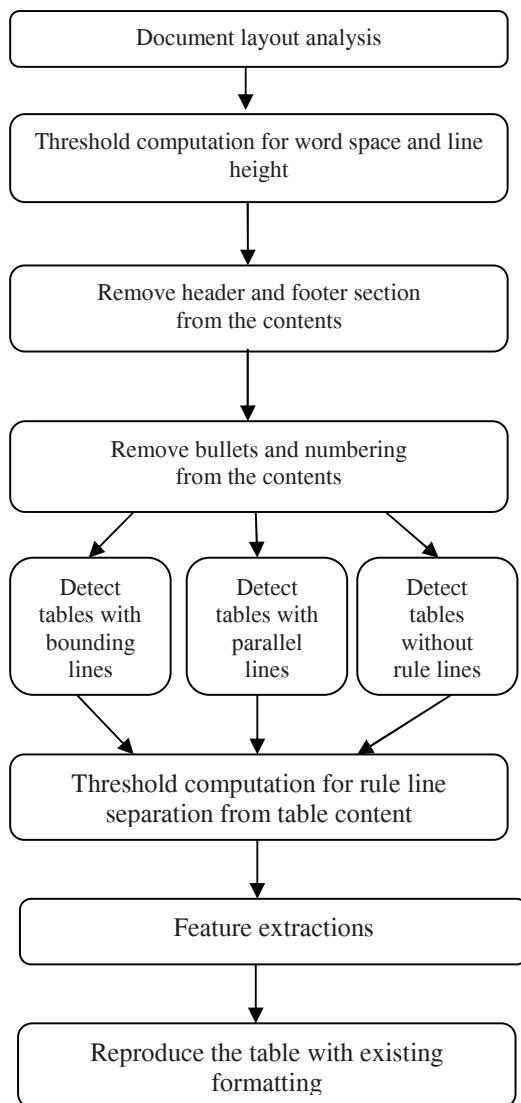


**Figure 2:** Pre-processing of document images

## Appendix B Microsoft SQL Server 2000, SQL Server 2005, and Oracle Database 10g Comparison

The following table shows the statements that are sent to the supported databases when database transactions are handled.

Microsoft Dynamics AX 4.0	SQL Server 2000	SQL Server 2005	Oracle Database 10g
First msbegin statement	SET TRANSACTION ISOLATION LEVEL READ COMMITTED	SET TRANSACTION ISOLATION LEVEL READ COMMITTED	No statement sent
First SQL DML statement inside a transaction scope	SET IMPLICIT_TRANSACTIONS ON	SET IMPLICIT_TRANSACTIONS ON	No statement sent
Final committing statement	COMMIT TRANSACTION	COMMIT TRANSACTION	COMMIT
Second msbegin statement	SET TRANSACTION I SOULATION LEVEL READ UNCOMMITTED	SET TRANSACTION I SOULATION LEVEL READ UNCOMMITTED	ROLLBACK
First SQL DML statement outside a transaction scope	SET TRANSACTION I SOULATION LEVEL READ COMMITTED	SET TRANSACTION I SOULATION LEVEL READ COMMITTED	ROLLBACK
select locked/locked	WITH (NOLOCK)	Not supported, so hint added	Not supported, so no hint added
Select optimized/concurrencyModel (ConcurrencyModel: Optimistic)	hint added to SELECT statement	No hint	No hint
Select pessimisticlock/concurrencyModel (ConcurrencyModel: Pessimistic)	WITH (UPDLOCK)	WITH (UPDLOCK)	FOR UPDATE OF clause added to SELECT statement



**Figure 3:** Flow of the algorithm

```

1. for x=1 to C do
2.   Y(x) ← Vertical projection profile
3.   T(x) ← Total of Y(x)
4.   R(x) ← T(x)/100
5.   if R(x) is greater than 5 then
6.     Partition(x) ← Y(x)
7.   end if
8. end for
Where
C - Number of column
T - Sum of vertical projection profile of a column
R - Threshold value for a text column
  
```

**Figure 4:** Algorithmic representation of layout analysis

the first word space in every line is compared with the maximum size of the word space of that particular line. When the first word space and the maximum size of the word space in a text line are equal, it is considered as a bulleted/numbered text line and eliminated from the content of the document.

#### **Threshold computation for word space and line height**

Threshold computation for word space is computed for a standard text line by eliminating other spaces greater than or less than the word space, using the assumption of having comparatively more number of gaps in a line (have a higher number of characters) that can produce a lengthy text line, which can be a standard text line as shown in Figure 5.

Initially, the maximum size of an inter-word gap in a standard text line is considered as a word space, and the height of the standard text line is considered as line height. Subsequently from the experiments, the range of thresholds for word space and line height are determined using the standard text line by

$$WS < ws < 2 * WS$$

$$LH < lh < 1.5 * LH$$

ws - computed word space from the standard text line

ws - threshold for word space

LH - computed line height from the standard text line

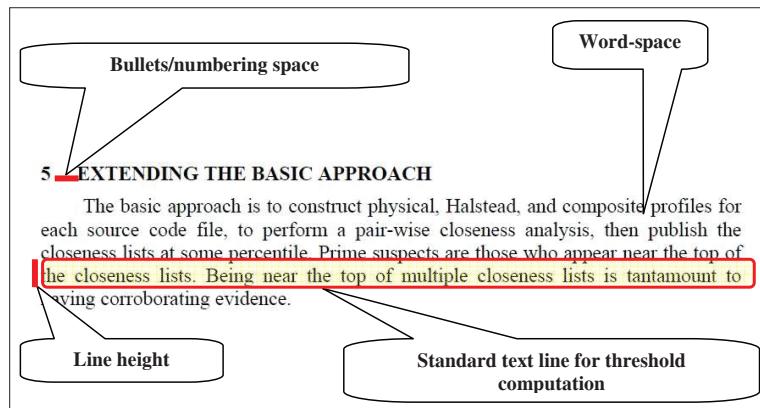
lh - threshold for line height

#### **Removing header and footer from the document**

Here we have compared the threshold for word space and line height with the consequent three lines at the top and bottom of the page. If the first line is a rule line and the other two lines are not related to tabulated data or paragraph content, then that part has been eliminated from the content page and considered as a header. Similarly, if the last line is a rule line and the other consecutive lines are not related to tabulated data or paragraph content, then that part is considered as the footer section and eliminated from the content.

#### **Detecting tables from the documents**

Part of the algorithm works on three consecutive passes such as extraction of the table with: i) fully and partially bounding lines; ii) parallel lines and iii) no ruling lines. Here we have checked the word space and the height of each line against the calculated thresholds and extracted



**Figure 5:** A sample document portion with text measurement

```

1. if  $S(1)$  is greater than  $ws$  and  $S(2)$  is less than or equal to  $ws$  or
2.  $S(1)$  equal to 0 and  $S(2)$  is less than or equal to  $ws$  and  $H(1)$  is less than  $lh$  or
3.  $S(1)$  is greater than  $ws$  and  $S(2)$  is less than or equal to  $ws$  and  $S(3)$  is less than or equal to  $ws$  then
4. start=2 else
5. if  $S(1)$  is greater than  $ws$  and  $S(2)$  is equal to 0 and  $S(3)$  is less than or equal to  $ws$  then
6. start=3 else
7. if  $S(1)$  equal to 0 and  $S(2)$  is greater than  $ws$  and  $S(3)$  is less than or equal to  $ws$  and  $H(1)$  is less than  $lh$  then
8. start=3 else
9. start=1
10. end if
11. end if
12. end if
13. if  $S(n)$  is greater than  $ws$  then
14. terminate=n-1 else
15. if  $S(n)$  is less than or equal to  $ws$  and  $S(n-1)$  is equal to 0 and  $S(n-2)$  is less than or equal to  $ws$  then
16. terminate=n-2 else
17. if  $S(n)$  is greater than  $ws$  and  $S(n-1)$  is equal 0 and  $S(n-2)$  is less than or equal to  $ws$  then
18. terminate=n-2 else
19. terminate=
20. end if
21. end if
22. end if

```

Where

$S$ -maximum length of word space in a line

$ws$ -threshold of a word space

start-start of the text content

terminate-end of the text content

$H$ -height of a line

**Figure 6:** Algorithmic representation for eliminating header-footer content

the text rows that belong to tabulated data. This part of the algorithm is similar to the one presented in the work by Jahan and Ragel (2014).

#### *Recognition of tables and reconstructing them with the formatting features*

From the three types of detected tables and their

contents, the tabular information is extracted using the fact that the rule lines can have more number of pixels than the tabulated contents in both vertical and horizontal direction.

The algorithm in Figure 7 shows the extraction of the thickness of the rule line. Similarly, all the other line features have been extracted.

```

1.  $X \leftarrow$ Horizontal projection profile
2.  $Size \leftarrow$ length of  $X$ 
3.  $maxTH \leftarrow$ maximum range of  $X$ 
4.  $minTH \leftarrow$ minimum of  $X$ 
5.  $midTH \leftarrow$ mid of  $X$ 
6. for  $n=1$  to number_of_lines do
7.   for  $i=1$  to size do
8.     if  $X(i)$  is not equal to  $minTH$  or  $midTH$  or  $maxTH$  and  $X(i+1)$  is equal to  $minTH$  or  $midTH$  or  $maxTH$  then
9.        $LS(n)=i$ 
10.    end if
11.    if  $X(i)$  is equal to  $minTH$  or  $midTH$  or  $maxTH$  and  $X(i+1)$  is not equal to  $minTH$  or  $midTH$  or  $maxTH$  then
12.       $LE(n)=i$ 
13.    end if
14.   end for
15.    $LT(n)=LE(n)-LS(n)$ 
16. end for

Where
LT-Thickness of the rule line
LS-Line start
LE-Line end

```

**Figure 7:** Algorithmic representation of part of feature extraction of a table

This recognition process executes some computational steps as given below:

(a) Threshold computation for separating lines from the tabular content and line identification

Threshold values were computed from a range of maximum number of pixels to separate the ruling lines with the content of the table using both horizontal and vertical projection profiles of the detected table. Then the identified lines were separated from the content.

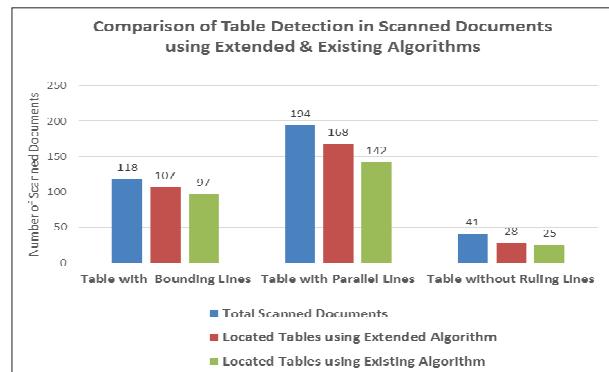
(b) Line feature extraction

Features of the tables' rule lines were extracted according to the following to reconstruct the table in the documents: i) the number of rows and columns; ii) the number of vertical and horizontal lines; iii) the height of each row; iv) the width of each column; v) the thickness of horizontal and vertical rule lines; and vi) the height and width of each table.

(c) Reproduce the tables with the formatting features

From the extracted line features, the tables were reconstructed by passing the extracted geometric features to construct an HTML file. Therefore, the users can edit and update the table quickly.

such as with bounding lines, parallel lines and without lines show comparatively higher number of table detection with the earlier algorithm. The tables with bounding lines show a comparatively higher detection. This type of tables can be extracted easily as it is bounded by at least a single vertical line. We compare the line heights with the thresholds to eliminate the interruption of main headings and equations within the text content. Therefore, they yield the highest document detection with the vertical ruled line information. For the tables with parallel lines, the current extended algorithm shows more improvements, which significantly increases the document detection rate. Although the new extended algorithm gives many improvements, there are several factors that still influence the accuracy of detection of tables in all three categories.



**Figure 8:** Comparison of table detection in scanned documents using extended and existing algorithms

	Right Index	Right Middle	Left Index	Left Middle	All Four Fingers
Time 1 (t1)	320	320	320	320	1280
Time 2 (t2)	320	320	320	320	1280
t1 + t2	640	640	640	640	2560

Table 2: Number of images used for performance evaluation.

the four fingers. This resulted in 1,280 ( $160 \times 4 \times 2$ ) fingerprint images per sensor. The subjects were requested to provide their fingerprint images again after a period of 6 weeks. At this second time instant, the above procedure was repeated to collect another 1,280 ( $160 \times 4 \times 2$ ) fingerprint images. So, our database consists of a total of 2,560 fingerprint images collected using Digital Biometric's optical sensor and an equal number of fingerprint images using the Veridion's solid-state optical sensor and the acquisition process was carried out by several volunteers who are research assistants in the Pattern Recognition and Image Processing Laboratory at Michigan State University. A live feedback of the acquired image was provided and the volunteers guided the subjects in placing their fingers in the center of the sensor and in an upright position. Due to this assistance provided to the subjects, most of the fingerprints were centered and there is no significant rotation in the acquired images.

The sensors were attached to IBM Thinkpads. Because of the time lag in providing the feedback to the user, some of the images were actually acquired at the time when the subject was in the process of removing his/her hand from the sensor. This resulted in some partial fingerprints. The volunteers observed that more subjects had greasy fingers during and after the lunch hour while more dry fingerprints were observed in the evening as compared to the morning session of the data acquisition process.

We collected another database to study the user adaptation issue in fingerprint image acquisition process. Eight subjects from our laboratory were requested to give one fingerprint sample every working day of the week for six consecutive weeks. The

15

	Right Index	Rigid Middle	Left Index	Left Middle	All Four Fingers
Thine 1 ([1])	330	330	330	330	1330
Thine 2 ([2])	320	330	330	330	1320
41 + 42	340	340	340	340	1340

(b)

A	A	A	A	A	A
A	A	A	A	A	A
A	A	A	A	A	A
A	A	A	A	A	A

(c)

**Figure 9A:** Reconstruction of the table: (a) table with bounded lines; (b) located table; (c) reconstructed table.

### Acknowledgments

We want to thank all the people who assisted us in making this book become a reality. The list of people is long—if we inadvertently missed anyone, we apologize. A special thanks goes to the following people on the Microsoft Dynamics product team:

- the following people on the Microsoft® SharePoint® product team:

  - Mette Nyberg, who worked with authors, their idiosyncrasies, and non-native English to get the book in a shape that would allow the Microsoft Press editors to actually start editing. She also made sure that the deadlines were kept within a safe horizon.
  - Hal Howard and Niels Bo Theilgaard, who sponsored the project.
  - The product team reviewers, who provided valuable feedback that made the book more accurate and a whole lot better.

Ajay Aggarwal	Jyoti Gavade
Hari Pulapaka	Peter Villadsen
Kenneth Puggard	Jens Klarlund Jensen
Lachlan Cuth	Nestor Serrano
Lauren Ricci	David Achter
Lei He	Ramana Parimi
Mark B. Madsen	Jeppe Øskar Meyer Larsen
Morten Jensen	Hua Chu
Per Vilksen	Ken Hirurup
Praveen Aradhanam	Anura Gunaratne
Sune Gynthersen	David Pokkuda
Uffe Kjall	Jens Møller-Pedersen

- We would also like to thank our external reviewers, who took time out of their busy schedules to add value to this book:  
Martin Fruergaard Laursen  
Anders Haage  
Oliver Morrison

Of course, we also want to thank the people at Microsoft Press who helped support us throughout the book writing and publishing process:

- Ben Ryan, who championed the book project at Microsoft Press.

<b>Musician/Critic</b>	<b>Album/Performance</b>
Amy Agler	Great American
Van Morrison	Irish Whiskies
Baroness Pugmire	Jesus Harkness Jesus
Letherette	Mobile Birth Machine
Lowdown Red	Donald Acreless
Li'l Ho	Human Peril
Mack B. Mabon	Joyful Order Major Lazer
Marcus Johnson	Mr. C.U.
Pit Michelin	Mr. Michelin
Soldado	Armenian
Soror Cythara	Donald Pendleton
Lilith St. Grail	Jesus Whales/Promised

(b)

xix

—(c)

**Figure 9B:** Reconstruction of the table: (a) table with parallel lines; (b) located table (c) reconstructed table

(a)

(b)

(c)

## 218 Part II Developing with Microsoft Dynamics AX

Table 9-1 AxBase Method Descriptions

Method	Description
updateNow	AxBase method that may be overridden by derived classes to implement document-specific updates (for example, posting the document or running some business logic). The method is called as the very last step of the updateNow method.
validateDocument	AxBase method that may be overridden by derived classes to implement document-wide business logic. This method is called immediately after processing an inbound document and right before the document is saved.

For improved productivity, the framework also implements operations by default. This means that, if you, for example, do not want to support `create`, you must override it in your derived document class and throw an exception explaining that the operation is not supported. The following example of an unsupported `create` action is from the actual implementation in the AxBase class, which is included with Dynamics AX:

```
public AxBatchKey create(AxTable _inTable,
    AxBatchActionPolicyInfo _actionPolicyInfo,
    AxTable _constraintList)
{
    ;
    throw error(strfmt("WS1504024", this.getName(), "create"));
}
```

When you create a completely new document, you must override the following methods:

- `getCreateList`
- `getConstraintList`
- `getAxbatch`

When you include more tables in the query, and the field values in one query rely on field values in, for example, the parent tables, you must also override the `prepareForSave` method. This is necessary only if you intend to support `create` or `readList` actions. In the following example, using the `prepareForSave` method of the `AxBatch` document class, the foreign key is populated with the values from the primary fields in the parent table. Pay particular attention to the lines in bold, which set up the relationships.

```
public boolean prepareForSave(AxStack axStack, AxTable str_dataSourceName)
{
    AxTable axtagTable = axTable("axTable");
    AxTable axtagTableInterval = axTable("axTableInterval");
    ;
    switch (classid == AxStack_top())
    {
        ;
    }
}
```

(d)

(e)

(f)

Table 4. Cross-year evaluation of text alignment software submissions for 2012 and 2013 w.r.t. precision. The darker a cell, the better the performance compared to the entire column.						
Software Submission		Observation Strategies of the 2012 Evaluation Corpus				Entire Corp.
Team	Year	Note	Random	Cyclic	Stratified	Entire Corp.
Gillman	2012	0.53606	0.97164	0.79933	(0.9999)	0.95151
Gillman	2013	0.25156	0.97124	0.77233	(0.9999)	0.95155
Baseline	2012	0.95354	0.98505	0.04444	0.98048	0.99744
Obersteiner	2012	0.86411	0.98501	0.87500	0.83904	0.96045
R. Terpstra	2012	0.48286	0.96042	0.87500	0.84553	0.97519
Kong	2012	0.53488	0.98267	0.75043	0.82079	0.90662
Kong	2013	0.50007	0.98266	0.75042	0.80790	0.90711
Koepnick	2013	0.80007	0.97166	0.71826	0.80790	0.96175
Schoenauer	2012	0.87103	0.98322	0.76462	0.82255	0.96319
Paliwal&Koepnick	2012	0.70408	0.68226	0.68478	0.82346	0.72824
Jayapal	2012	0.86978	0.57056	0.24240	0.74551	0.88731
Jayapal	2013	0.87320	0.44854	0.17098	0.72738	0.26708
Sánchez-Vega	2012	0.95164	0.98576	0.11664	0.95980	0.84231
Sánchez-Vega	2013	0.95164	0.98576	0.11664	0.95980	0.57511
Software Submission		Observation Strategies of the 2013 Evaluation Corpus				Entire Corp.
Team	Year	Note	Random	Cyclic	Stratified	Entire Corp.
Gillman	2012	0.95492	0.98542	0.95965	0.95950	0.92209
Jayapal	2012	0.85942	0.95842	0.95965	0.95950	0.84507
Baseline	2013	0.87041	0.96101	0.97835	0.91167	0.95147
R. Terpstra	2013	0.90906	0.96996	0.95014	0.95070	0.95048
Dietmar	2013	0.90906	0.96996	0.95014	0.95070	0.95048
Gillman	2013	0.81826	0.95972	0.97273	0.95991	0.95991
Jayapal	2013	0.88003	0.99984	0.97773	0.99950	0.88487
Schoenauer	2013	0.89003	0.92335	0.98008	0.94555	0.95631
Koepnick	2012	0.82675	0.95942	0.95950	0.95950	0.95665
Kong	2012	0.32675	0.91810	0.95950	0.95950	0.95665
Kong	2013	0.87076	0.98967	0.85421	0.85399	0.85297
Schoenauer	2012	0.61678	0.97581	0.81531	0.84748	0.84577
Kong	2013	0.87076	0.98967	0.85421	0.85399	0.85297
R. Terpstra	2012	0.81313	0.93881	0.81159	0.95666	0.81313
Paliwal&Koepnick	2012	0.90700	0.95137	0.81234	0.97476	0.90700
Schoenauer	2013	0.60932	0.92073	0.64994	0.70708	0.75144

#### **Appendix 3: U.S. Census Bureau 2011 “People Quick Facts”**

<b>People Quick Facts</b>	<b>United States</b>
Population, 2011 estimate	311,591,917
Population, 2010 (April 1) estimates base	308,745,538
Population, percent change, April 1, 2010 to July 1, 2011	0.9%
Population, 2010	308,745,538
Persons under 5 years, percent, 2011	6.3%
Persons under 18 years, percent, 2011	23.7%
Persons 65 years and over, percent, 2011	13.3%
Female persons, percent, 2011	50.8%
White persons, percent, 2011 (a)	78.1%
Black persons, percent, 2011 (a)	13.1%
American Indian and Alaska Native persons, percent, 2011 (a)	1.2%
Asian persons, percent, 2011 (a)	5.0%
Native Hawaiian and Other Pacific Islander persons, percent, 2011 (a)	0.2%
Persons reporting two or more races, percent, 2011	2.3%
Persons of Hispanic or Latino Origin, percent, 2011 (b)	16.7%
White persons not Hispanic, percent, 2011	63.4%

Hospital Patient Number
Patient Social Security Number
Patient Racial Background
Patient Birth Date
Patient Gender
Visit Date
Discharge Date
Principal Diagnosis Code (ICD9)
Procedure Codes (up to 14)
Medication Codes (up to 14)
Primary Physician ID Number
Type of Physician
Total Medication Charges

56 | CHAPTER

**Access the Settings Manager for your Flash Player**

- 1 Open an application in Flash Player.
- 2 Right-click and select Settings.  
The Adobe Flash Player Settings dialog box appears.
- 3 Select the Privacy tab (on the far left).
- 4 Click the Advanced button.

Flash Player launches a new browser window and loads the

control so that the team member can make all kinds of changes on the site (including adding and removing user accounts), you need to know how to modify permissions. Here's how to do that:

1. On the Admin Overview page, click Manage Team Sites in the Team Sites And Documents category.
  2. On the Team Site Settings page, click Site Permissions in the Users And Permissions area.
  3. Click Grant Users Permission Directly (shown in Figure 3-16), and then click the permission level you want to assign to the team member.
  4. Type the text for the email message if you want to send one, and click OK. The team member is assigned the permission level you selected.

#### **Other resources**

(g)

(h)

(i)

**Figure 10:** Factors contributing to false detection of table: (a) multiple text lines within a column in a particular row; (b) heterogeneous font types and styles on a single page; (c) different font types and styles; (d) different inter-word space; (e) grey and black table cell; (f) full lengthy table; (g) unorganised table structure; (h) numbering style 1; (i) numbering style 2.

The following section shows some examples of located and reconstructed tables from the printed documents. Figure 9 depicts how a reconstructable table is extracted from a scanned document. Initially, the table portion is located on the document page, and the extracted information from the table is reused to reconstruct the table without tabulated content. We consider that recognition and producing the table content is out of the scope of this paper. However, such recognition can be performed with well-established techniques like OCR.

As indicated in the results, we have identified several factors that influence table detection in scanned documents. Among the 41 collected documents that contain tables without ruling lines, 13 tables have not been located in the documents.

*Problem 1:* The table has multiple text lines within one column and a single line in the other column in a particular row as shown in Figure 10 (a). As a result, when comparing the threshold value with inter-column space, the tabulated text lines have not been extracted except the first text line. Here, five of the documents have not been detected.

*Problem 2:* The document contains tables and their contents only on the entire pages that lead to the deficiency of text lines for threshold computation. Therefore, threshold computation is proceeded using the tabulated data rather than the standard text line that leads to erroneous detection. Here, two of the documents have not been identified.

*Problem 3:* The differences in font type and style within a single page as shown in Figure 10 (b) leads to an erroneous situation when computing a standard word space. Two of the documents could not be identified.

*Problem 4:* The optimisation process in threshold calculation is also one of the reasons for the erroneous situation in four of the documents. In addition, as this type of a table solely depends on the word space threshold value rather than the ruling lines of tables, it shows a challenge to produce a higher accuracy.

In the tables with parallel lines, among 194 documents, 26 documents have not been identified. Many reasons obstruct the detection of these tables using the given algorithm.

*Problem 1:* The table, as shown in Figure 10(e), comprises a black and grey cell structure that completely interrupt the threshold computation as our algorithm only works on grey and white background, and four documents produced error results.

*Problem 2:* In some cases the table is not correctly structured as shown in Figure 10(g), which leads to the algorithm being more competitive when differentiating the table component from the text elements in the document. Therefore, one of the documents could not produce the positive output for table detection.

*Problem 3:* In some cases, the document contains only a large table and the contents of a minimum number of extra text lines on the particular page. Such documents can lead to the deficiency of enough text lines for threshold computation. This type of tables are not detected in two of the documents.

*Problem 4:* The given algorithm has the facility to eliminate bullets and numbering from the text contents as shown in Figure 10(h). However, the space existing in the different numbering structure as illustrated in figure 10(i) interrupts the threshold calculation and comparison. Therefore, these types of tables are not detected in three of the documents.

*Problem 5:* In some documents, the table has multiple text lines within one column and a single line in the other column in a particular row as shown in Figure 10(a). As a result, when comparing the threshold value with inter-column space, the tabulated text lines have not been extracted except the first text line. Here, the tables have not been detected in eleven of the documents.

*Problem 6:* Sometimes the document can have different font types and styles. For example, a paper that has a regular text and programming codes can be differentiated by various font types, styles and inter-word gaps as shown in Figures 10(c) and 10(d). This will lead to a problem in threshold calculation as it has different word-spaces compared to the regular text line. This has caused an adverse effect on the detection of five of the document.

From the analysis of the reasons for errors, most of the individual components of the algorithm successfully yielded the output except a few steps of the algorithm. As we have analysed the results for 194 tables with parallel lines, the performance of individual components of the algorithm is as follows.

- Component 1 - almost all the documents completely give the precise output.
- Component 2 - three of the documents yield an error.
- Component 3 - eleven documents show error.
- Component 4 - almost all the documents work completely with this component
- Component 5 - eleven documents show error.
- Component 6 - almost all the detected documents produce a positive output.

Similarly, as we have analysed the results for 41 tables without ruling lines, the performance of the individual component of the algorithm is as follows.

- Component 1 - almost all the documents completely give the precise output.
- Component 2 - almost all the documents completely give the precise output.
- Component 3 - thirteen documents do not show precise output.
- Components 4, 5 and 6 - all the remaining documents produce positive results.

Time complexity is one of the measurements that plays a significant role when we select an efficient algorithm for a problem domain. Here, the time complexity of the proposed algorithm is linear and produces  $O(N)$  complexity, where  $N$  is the number of text or rule lines on a document page. Initially, in the first component of the algorithm-threshold computation for layout analysis, the complexity is  $O(1)$  as it only handles the vertical projection profile of the document image. The second and third step of the algorithm such as computing the number of spaces in each line, handling bullets and numbering lines and calculating threshold values for line height and word space, completely depend on the number of visible rule or text lines on the document page, and they still behave as  $O(N)$ -linear complexity. The next subsequent step of removing the header and footer is in  $O(1)$  complexity as it follows a single line of the document image. The other major stage of the algorithm locating of table behaves linearly as it still depends on the number of the entire rule or text lines of the document page and the complexity is  $O(N)$ . The final stage of the algorithm recognising and reproducing the table only depends on the pixels of a few tabular lines. This table portion is comparatively smaller in size than the text portion of a document page and can be negligible. Therefore, the algorithm completely depends on the number of pixels and traverse in the vertical direction of the document page. The computational steps are still determined by the number of lines on a document page. Hence, from the overall analysis of the complete algorithm, it behaves as linear, and the time complexity of the complete algorithm is  $O(N)$ . The algorithm works on a single document page at a time, and the size of a document image may very much depend on the document.

## CONCLUSION

In this paper, we have proposed and evaluated an extended new algorithm for locating tables from scanned documents. From the experiments, the system

has shown comparatively much higher detection ability than the existing algorithm addressed in Jahan and Ragel (2014). We have determined an automated local threshold for line height and word space, which can be determined from a particular document and completely depend on it. Although the intermediate processes such as i) layout analysis using local threshold; ii) eliminating bullets and numbering lines and iii) eliminating header and footer section from the text content results in more improvement within the extending algorithm, the process of threshold computation for word space and line height can be enhanced in the future to achieve robustness of the system. The motivation towards the extraction of tabular features can be used to reproduce the detected tables with their real formatting features for future editing and updating purposes. The algorithm has a linear complexity of  $O(N)$ , which will be advantageous while processing a large number of scanned documents during the rapid conversion of the traditional library system into a digital system.

---

## REFERENCES

1. Ajward S., Jayasundara N., Madushika S. & Ragel R. (2010). Converting printed Sinhala documents to formatted editable text. *Proceedings of the 5<sup>th</sup> International Conference on Information and Automation for Sustainability (ICIAFs)*, volume 5, Colombo, Sri Lanka, 17 – 19 December, pp. 138 – 143.  
DOI: <http://dx.doi.org/10.1109/iciafs.2010.5715649>
2. Chen J. & Lopresti D. (2011). Table detection in noisy off-line handwritten documents. *Proceedings of the 11<sup>th</sup> International Conference on Document Analysis and Recognition*, Beijing, China, 18 – 21 September, pp. 399 – 403.
3. Cesarini F., Marinai S., Sarti L. & Soda G. (2002). Trainable table location in document images. *Proceedings of the 16<sup>th</sup> International Conference on Pattern Recognition*, Quebec City, Canada, 11 – 15 August, pp. 236 – 240.  
DOI: <http://dx.doi.org/10.1109/icpr.2002.1047838>
4. Gatos B., Danatsas D., Pratikakis I. & Perantonis S.J. (2005). Automatic table detection in document images. *Proceedings of the 3<sup>rd</sup> International Conference on Advances in Pattern Recognition*, Bath, UK, 22 – 25 August, pp. 609 – 618.  
DOI: [http://dx.doi.org/10.1007/11551188\\_67](http://dx.doi.org/10.1007/11551188_67)
5. Gatos B., Pratikakis I. & Perantonis S.J. (2004). An adaptive binarization technique for low-quality historical documents. *Proceedings of the 6<sup>th</sup> International Workshop on Document Analysis Systems*, Florence, Italy, 8 – 10 September, pp. 102 – 113.  
DOI: [http://dx.doi.org/10.1007/978-3-540-28640-0\\_10](http://dx.doi.org/10.1007/978-3-540-28640-0_10)
6. Harit G. & Bansal A. (2012). Table detection in document images using header and trailer patterns. *8<sup>th</sup> Indian Conference on Computer Vision, Graphics and Image Processing*, Mumbai, India, 16 – 19 December, (Pages?)

7. Hu J., Kashi R., Lopresti D. & Wilfong G. (2000). Medium-independent table detection. *SPIE Document Recognition and Retrieval VII*: 291 – 302.
8. Jahan M.A.C.A. & Ragel R.G. (2014). Locating tables in scanned documents for reconstructing and republishing. *Proceedings of 7<sup>th</sup> International Conference on Information and Automation for Sustainability*, Colombo, Sri Lanka, 20 – 22 December, (pages?).
9. Fang J., Gao L., Bai K., Qiu R., Tao X. & Tang Z. (2011). A table detection method for multipage PDF documents via visual separators and tabular structures. *Proceedings of the 11<sup>th</sup> International Conference on Document Analysis and Recognition (ICDAR'11)*, Beijing, China, 18 – 21 September, pp. 779 – 783.
10. Kasar T., Barlas P., Adam S., Chetalain C. & Pacquet T. (2013). Learning to detect tables in scanned document images using line information. *Proceedings of the 12<sup>th</sup> International Conference on Document Analysis and Recognition*, Washington, DC, USA, 25 – 28 August, (Pages?)
11. Kieninger T.G. (1998). Table structure recognition based on robust block segmentation. *Proceedings of the SPIE – The International Society for Optical Engineering*, volume 3305, pp. 22 – 32.
12. Kieninger T. & Dengel A. (2001). Applying the T-RECS table recognition system to the business letter domain. *Proceedings of the 11<sup>th</sup> International Conference on Document Analysis and Recognition*, Beijing, China, 18 – 21 September, pp. 518 – 522.  
DOI: <http://dx.doi.org/10.1109/ICDAR.2001.953843>
13. Laurentini T. & Vaida P. (1992). Identifying and understanding tabular material in compound documents. *Proceedings of the International Conference on Pattern Recognition*. pp. 405 – 409. (venue? Date?)  
DOI: <http://dx.doi.org/10.1109/icpr.1992.201803>
14. Mandal S., Chowdhury S.P., Das A.K. & Chanda B. (2006). A simple and effective table detection system from document images. *International Journal on Document Analysis and Recognition* 8(2): 172 – 182.
15. Namboodiri A. (2004). On-line handwritten document understanding. *PhD thesis*, Michigan State University, USA.
16. Oro E. & Ruffolo M. (2009). PDF-TREX: an approach for recognizing and extracting tables from PDF documents. *Proceedings of the 10<sup>th</sup> International Conference on Document Analysis and Recognition*, Barcelona, Spain, 26 – 29 July, pp. 906 – 910.
17. Shafait F. & Smith R. (2010). Table detection in heterogeneous documents. *Proceedings of the 9<sup>th</sup> IAPR International Workshop on Document Analysis Systems*, Boston, MA, USA , 9 – 11 June, pp. 65 – 72.  
DOI: <http://dx.doi.org/10.1145/1815330.1815339>
18. Smith R. (2009). Hybrid page layout analysis via tab-stop detection. *Proceedings of the 10<sup>th</sup> International Conference on Document Analysis and Recognition*, Barcelona, Spain, 26 – 29 July, pp. 241 – 245.  
DOI: <http://dx.doi.org/10.1109/icdar.2009.257>
19. Wang Y., Phillips I.T. & Haralick R. (2001). Automatic table ground truth generation and a background-analysis-based table structure extraction method. *Proceedings of the 6<sup>th</sup> International Conference on Document Analysis and Recognition*, Washington, USA, 10 – 13 September, pp. 528 – 532.
20. Watanabe T., Naruse H., Luo Q. & Sugie S. (1991). Structure analysis of table-form documents on the basis of the recognition of the vertical and horizontal line segments. *Proceedings of the 1<sup>st</sup> International Conference on Document Analysis and Recognition*, Saint – Malo, France, 30 September – 02 October, pp. 638 – 646.
21. Wonkyo S., Hyung I.K. & Nam I.C. (2015). Junction-based Table Detection in camera captured document images. *International Journal of Document Analysis and Recognition* 18(1): 47 – 57.
22. Yildiz B., Kaiser K. & Miksch S. (2005). Pdf2table: a method to extract table information from pdf files. *Proceedings of the 2<sup>nd</sup> Indian International Conference on Artificial Intelligence (IICAI'05)*, Pune, India, 20 – 22 December, pp. 1773 – 1775.
23. Zanibbi R., Blostein D. & Cordy J.R. (2003). A survey of table recognition: models, observations, transformations, and inferences. *International Journal of Document Analysis and Recognition* 7: 1 – 16.