

Document Layout Analysis for Semantic Information Extraction

W. T. Adrian^{1,2}, N. Leone¹, M. Manna¹, C. Marte¹

¹Department of Mathematics and Computer Science, University of Calabria, Italy

²AGH University of Science and Technology, Krakow, Poland

{w.adrian,leone,manna,marte}@mat.unical.it



AI*IA 2017



UNIVERSITÀ DELLA CALABRIA

BACKGROUND

- Due to the increasing volume of electronic documents, it is important to devise and refine techniques to automatically analyze and understand their content
- Document Layout Analysis (DLA)** is the process of decomposing a document into its component regions and understanding their functional roles and relationships
- The field of DLA is concerned with geometrical and logical labeling of document content
- Among the phases of DLA, *page decomposition* is a key phase that aims to identify maximal “homogeneous” regions of the document
- Existing approaches to page decomposition exploit different approaches/techniques:
 - Bottom-up* (or aggregating method)
 - Top-down* (or divisive method)

MOTIVATION AND OBJECTIVES

- Despite the progress observed over the last years, **challenges are still open**: they range from accurately detecting content boxes to classifying them into semantically meaningful classes
- The variety of document styles and formats makes non-trivial the task of automatically recognize the structure of arbitrary documents
- Our interest in document layout analysis mainly derives from actual needs of **improving Information Extraction** performances over complex real-world documents sharing some “key information”
- Curriculum Vitae represent one possible interesting use-case

DOCUMENT CHARACTERIZATION

- From a logical viewpoint each word can be characterized as a tuple $w = (id, val, x_1, y_1, x_2, y_2, \kappa)$

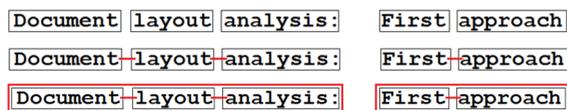


- Definition 1.** Give to words w and w' , we say that w' follows w if $x_2 < x'_1$ and if at least one of the following conditions is satisfied:

- $y'_2 \leq y_1 \leq y'_1 \wedge x_1 - x_1 < \kappa$;
- $y'_2 \leq y_2 \leq y'_1 \wedge x_1 - x_2 < \kappa$;
- $y_2 < y'_2 < y'_1 < y_1 \wedge x_1 - x_2 < \kappa$;

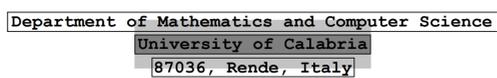


- Definition 2.** A **line block** is formally defined as any maximal connected component of the graph associated to the document. By definition, a word that is not followed by and does not follows any other word is trivially a line block.

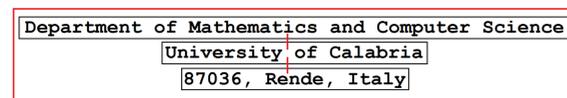


- Definition 3.** Give two line blocks LB_{Down} and LB'_{Up} , we say that LB_{Down} and LB'_{Up} are **linked lines** if $y_T < y'_B$ and if at least one of the following conditions is satisfied:

- $x'_L \leq x_L \leq x'_R$;
- $x'_L \leq x_R \leq x'_R$;
- $x_L < x'_L < x'_R < x_R$.



- Definition 4.** A **text block** is formally defined as any maximal connected component of the graph associated to line blocks. Moreover, a text block that is not connected to any other line block is, trivially, a text block.



Real result on a PDF document

LOGIC-BASED ENCODING

- Our input consists of a set of words that can be represent by a fact of the form $word(Id, Val, X1, Y1, X2, Y2, K)$
- To identify pairs of words connected by the relation *follows* we use a logic rule for each condition e.g.:

```
follows(IdL, IdR) :- word(IdL, ValL, X1L, Y1L, X2L, Y2L, K1),
                    word(IdR, ValR, X1R, Y1R, X2R, Y2R, K2),
                    Y2R <= Y1L, Y1L <= Y1R, X2L < X1R, Delta = X1R - X2L, Delta <= K1.
```
- To identify the first element of a text line, we use a concept of an “ancestor”, which is inferred from the *follows* relation, and then keep it in *inLine* predicate:

```
linkedWord(IdL) :- follows(IdL, IdR). linkedWord(IdR) :- follows(IdL, IdR).
isolateWord(Id) :- word(Id, _, X1, Y1, X2, Y2, K), not linkedWord(Id).
ancestorOf(IdA, IdB) :- follows(IdA, IdB), IdA < IdB.
ancestorOf(IdA, IdC) :- ancestorOf(IdA, IdB), follows(IdB, IdC), IdA < IdC.
child(Id) :- ancestorOf(_, Id). inLine(Id, Id) :- isolateWord(Id).
inLine(IdA, IdA) :- ancestorOf(IdA, IdB), not child(IdA).
inLine(IdA, IdB) :- ancestorOf(IdA, IdB), not child(IdA).
```

- Afterwards, we calculate the bounding box of the text line:

```
blockLine(IDln, X1, Y1, X2B, Y2B) :- word(IDln, _, X1, Y1, X2, Y2, _),
                                   word(IDB, _, X1B, Y1B, X2B, Y2B, _),
                                   inLine(IDln, IDB), not overcome(IDln, X2B).
```
- Similarly, we build rules that encode conditions for the linked lines.

- Finally, we find the bounding box of the entire text block:

```
ascissaTL(IDblk, X1) :- inBlock(IDblk, IDln), blockLine(IDln, X1, Y1, X2, Y2),
                       not existsLowerThanTL(IDblk, X1).
ordinataTL(IDblk, Y1) :- inBlock(IDblk, IDln), blockLine(IDln, X1, Y1, X2, Y2),
                        not overcomeTL(IDblk, Y1).
angleTopLeft(IDblk, X1, Y1) :- ascissaTL(IDblk, X1), ordinataTL(IDblk, Y1).
```

CONTENT RECOGNITION

- We divide the document in *homogeneous section* and proceed with the structure recognition technique (Fig.1)
- We use the external tool Quablo (<http://quablo.eu>) to recognize a table or *draw* one from fixed section of the PDF document, obtaining a grid representation (Fig.2)
- Consecutively, we normalize the grid representation with the use of a *heading ontology*, in order to obtain a **two-dimensional representation** of the document (Fig.3)

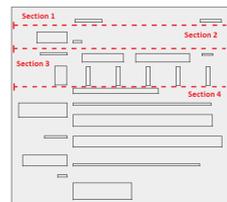


Figure 1

Figure 2

Figure 3

- This is the base to perform Information Extraction with *extraction rules* that use the spatial context of information

CONTRIBUTION

- We use **Logic Programming** to provide a declarative characterization of the basic geometric properties of text units/block
- We split the document in sections, each containing some “homogeneous” layout
- We exploit external tools for table recognition to detect tables
- We use **Semantic Annotation** to identify/classify atomic text units
- We reason over geometric and semantic information to identify blocks that are “logically connected”

RESEARCH PROPOSAL – FUTURE WORK

Geometric and Semantic Aspects for Document Layout Analysis

Cinzia Marte

For the near future we would like to:

- dynamically* refine and improve the parameter κ deduced in the developed approach, used to provide an upper-bound to the width of a white space placed between consecutive words
- increase the graph representing words and their connection by adding values on the edges, estimating these values both from syntactic and semantic point of view. In particular, we think to define a multi-variable function (one variable for each aspect considered, as semantical, syntactic, etc.), that represents the probability that two words are consecutive
- determine a reading flow* in order to merge some areas above detected and *automatically* deduce homogeneous section of the document
- use *Disjunctive Logical Programming* to refine the characterization of fixed part of the document and to provide an improved declarative representation of the basic logical properties of the document
- combine efficiently this information, exploiting as much semantic information as possible, in order to obtain a truthful/faithful page segmentation, useful to get better Information Extraction performances over document as Curriculum Vitae