

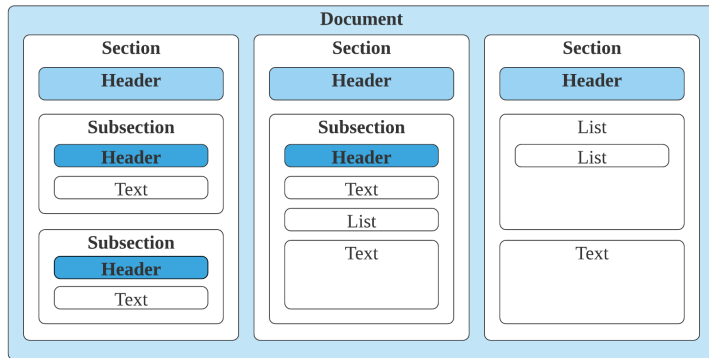
Извлечение иерархической логической структуры из текстовых документов в формате docx

Богатенкова Анастасия

5 ноября 2020 г.

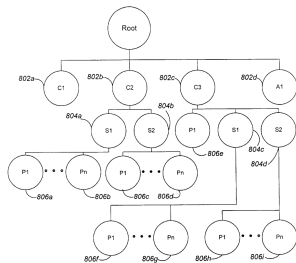


Многие документы имеют логическую структуру, выделение которой может помочь при решении задач автоматизированного анализа документов.

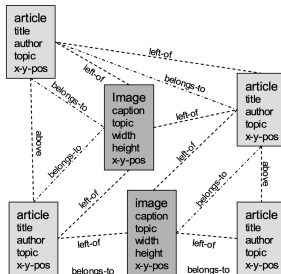


- В настоящее время большое количество документов создаётся и хранится в формате docx.
- Однако данный формат позволяет описывать только физическую структуру документа, то есть описывается то, как **выглядит документ**.
- Хотелось бы уметь выделять из подобных документов логическую структуру

Существуют различные способы представления структуры документа. Мы будем извлекать иерархическую структуру в виде дерева, так как многие документы состоят из последовательности вложенных друг в друга частей.



Tree structure



Graph structure

(a) $\begin{aligned} &[0] \langle \text{START} \rangle \rightarrow \langle \text{TITLE} \rangle \langle \text{COLUMN} \rangle \langle \text{COLUMN} \rangle \\ &[0] \langle \text{START} \rangle \rightarrow \langle \text{TITLE} \rangle \langle \text{COLUMN} \rangle \\ &[1] \langle \text{TITLE} \rangle \rightarrow \langle \text{text.line} \rangle \\ &[1] \langle \text{COLUMN} \rangle \rightarrow \langle \text{TEXT_BLOCKS} \rangle \\ &[0] \langle \text{TEXT_BLOCKS} \rangle \rightarrow \langle \text{TEXT_BLOCK} \rangle \langle \text{space} \rangle \langle \text{TEXT_BLOCKS} \rangle \\ &[0] \langle \text{TEXT_BLOCKS} \rangle \rightarrow \langle \text{TEXT_BLOCK} \rangle \\ &[1] \langle \text{TEXT_BLOCK} \rangle \rightarrow \langle \text{TEXT_LINES} \rangle \\ &[0] \langle \text{TEXT_LINES} \rangle \rightarrow \langle \text{text.line} \rangle \langle \text{newline} \rangle \langle \text{TEXT_LINES} \rangle \\ &[0] \langle \text{TEXT_LINES} \rangle \rightarrow \langle \text{text.line} \rangle \end{aligned}$

(b) $\begin{aligned} &\langle \text{START} \rangle \rightarrow \langle \text{TITLE} \rangle \langle \text{COLUMN} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{COLUMN} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{TEXT_BLOCKS} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{TEXT_BLOCK} \rangle \langle \text{space} \rangle \langle \text{TEXT_BLOCKS} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{text.line} \rangle \langle \text{newline} \rangle \langle \text{TEXT_BLOCK} \rangle \langle \text{space} \rangle \langle \text{TEXT_BLOCKS} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{text.line} \rangle \langle \text{newline} \rangle \langle \text{text.line} \rangle \langle \text{space} \rangle \langle \text{TEXT_BLOCKS} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{text.line} \rangle \langle \text{newline} \rangle \langle \text{text.line} \rangle \langle \text{space} \rangle \langle \text{TEXT_BLOCK} \rangle \\ &\rightarrow \langle \text{text.line} \rangle \langle \text{text.line} \rangle \langle \text{newline} \rangle \langle \text{text.line} \rangle \langle \text{space} \rangle \langle \text{text.line} \rangle \end{aligned}$

Fig. 6. Representing a document in terms of formal grammars: (a) example of a stochastic context free grammar that derives a text document with a title. Upper case symbols refer to non-terminal symbols, while lower case symbols show terminal symbols. (b) a sample document with a title and three lines, derived using the grammar in (a)

Formal grammar
representation

- Провести обзор некоторых форматов документов;
- Провести обзор способов представления логической структуры документа;
- Описать особенности формата docx;
- Описать структуру, которую необходимо извлечь;
- Реализовать метод извлечения структуры и провести экспериментальную проверку реализованного метода;
- Провести оценку качества.