

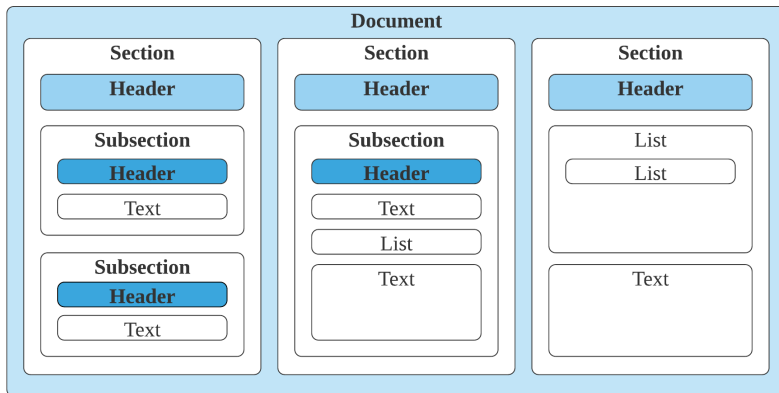
Извлечение иерархической логической структуры из текстовых документов в формате docx

Богатенкова Анастасия

5 ноября 2020 г.

Актуальность темы

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



Почему docx?

- ▶ 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

План работы

- ▶ 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.