



Санкт-Петербургский государственный университет
Кафедра системного программирования

Автоматическая подстановка параметров текстовых макросов в корректной форме

Карими Хурматулла, Группа 22.M07-мм

Научный руководитель: к.ф.-м.н. Д.В.Лутцев, доцент Кафедры системного программирования

Санкт-Петербург
2022

- Веб сайты, использующие язык разметки
 - ① Википедия
 - ② Github
 - ③ ReStructuredText
- Для английского языка подстановка простая
- Для русского подстановка требует менять форму словосочетания

Разработать и внедрить систему автоматической замены терминов в тексте в требуемой форме на русском языке.

- Изучить инструменты анализа текстов на естественных языках, позволяющие анализировать предложения.
- Выбрать модели машинного обучения, которые позволят, обучившись на корректных текстах выбирать нужные формы для подстановки фрагментов в текст.
- Провести эксперименты с этими моделями, выбрать наиболее подходящую.
- Реализовать прототип инструмента, выполняющего макроподстановку словосочетаний в корректном падеже и числе.

Инструменты, которые мы рассмотрели для анализа текстов:

- NLTK (Natural language toolkit).
- Py morphology2
- LemmInflect

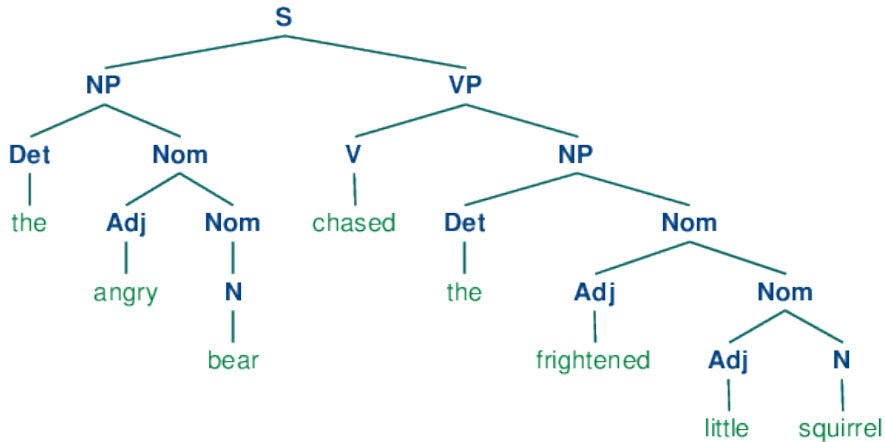
NLTK (Инструментарий естественного языка)¹

Это платформа, которая создает программы на Python, работающие с данными на человеческом языке. NLTK содержит библиотеки классификации текста, кластеризации и обработки текста для обозначения, синтаксического анализа, классификации и создания частей речи из наших данных. NLTK слишком велик, чтобы объяснять, но мы сосредоточимся на очень специфической области нашей работы:

- Рекурсивная структура для структуры грамматики.
- Рекурсивный-нисходящий-синтаксический анализ для разбора структурированной грамматики

¹<https://www.nltk.org/>

Рекурсивная структура



Рекурсивный-нисходящий-синтаксический анализ

1. Initial stage

S

the dog saw a man in the park

2. Second production



the dog saw a man in the park

3. Matching *the*



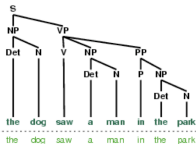
the dog saw a man in the park

4. Cannot match *man*



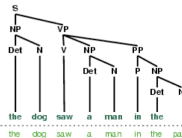
the dog saw a man in the park

5. Completed parse



the dog saw a man in the park

6. Backtracking



the dog saw a man in the park

Руморphy2 это морфологический анализатор, который анализирует русские тексты с помощью словаря орпесодрога. Алгоритм Руморphy выполняет морфологическую обработку на основе грамматической характеристики типа (слова, лемматизация), но если слово не существует в словаре, поэтому предиктор в руморphy2 объединит два алгоритма:

- 1 По префиксу
- 2 В конце слов

В то же время, когда мы анализируем слова, мы столкнемся с несколькими состояниями слов и оценками, затем, выбрав состояние слова с наивысшим баллом Мы можем исправить предложение

LemmInflect использует словарный подход для лемматизации английских слов и преобразования их в формы, заданные предоставленным пользователем тегом Universal Dependencies или Penn Treebank. Библиотека работает со словами, не входящими в словарный запас (OOV), применяя методы нейронных сетей для классификации словоформ и выбора соответствующих правил морфинга.

- Выбрать модели машинного обучения, которые позволят, обучившись на корректных текстах выбирать нужные формы для подстановки фрагментов в текст.
- Провести эксперименты с этими моделями, выбрать наиболее подходящую.
- Реализовать прототип инструмента, выполняющего макроподстановку словосочетаний в корректном падеже и числе.

- Bird, Steven. "NLTK: the natural language toolkit." Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions. 2006.
- <https://lemminflect.readthedocs.io/en/latest/>
- <https://pymorphy2.readthedocs.io/en/stable/>