

**МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ
(НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ)
Факультет информационных технологий и прикладной математики
Кафедра вычислительной математики и программирования**

**Отчет по лабораторной работе №4
«Построение сниппетов»
по курсу
«Обработка текстов на естественном языке»**

Группа: М80-108М-17
Выполнил: Забарин Н.И.
Преподаватель: А.Л. Калинин

Москва, 2018

Задание

Необходимо добавить в поисковую систему построение цитат (сниппетов), реферирование документов, найденных по запросу.

Сниппеты должны содержать слова запроса и давать пользователю представление о том, насколько документ отвечает поисковому запросу. Длина сниппета должна быть ограничена двумя-тремя строчками.

В отчете нужно привести описание алгоритма построения сниппетов, примеры.

Решение

Генерировать сниппеты было решено динамически, так как эта группа алгоритмов лучше отражает причины по которым документ попал в выборку.

Для реализации было внесено несколько исправлений в итераторы, теперь вместе с найденными документом возвращаются и слова запроса которые присутствуют в документе. Используя эту информацию и текст статьи будем искать окно с ограниченной длиной и максимальным числом слов из запроса. Код ищущий сниппет:

```
def extract_snippet(text, q):
    l, r = 0, 1
    curlen = 0
    curans = 0
    spt = []
    term = ""

    ans = []
    mx = 0

    def check_term(term):
        for i in q:
            if term.startswith(i):
                return True
        return False

    i = 0
    while i < len(text):
        if text[i] in stop_sym:
            while text[i] in stop_sym:
                term += text[i]
                i += 1

            curlen += len(term)
            spt.append(term)
            if curlen > SNIPPET_MAX_LEN:
                if check_term(spt[0]):
                    curans -= 1
                spt = spt[1:]
            if check_term(term):
                curans += 1

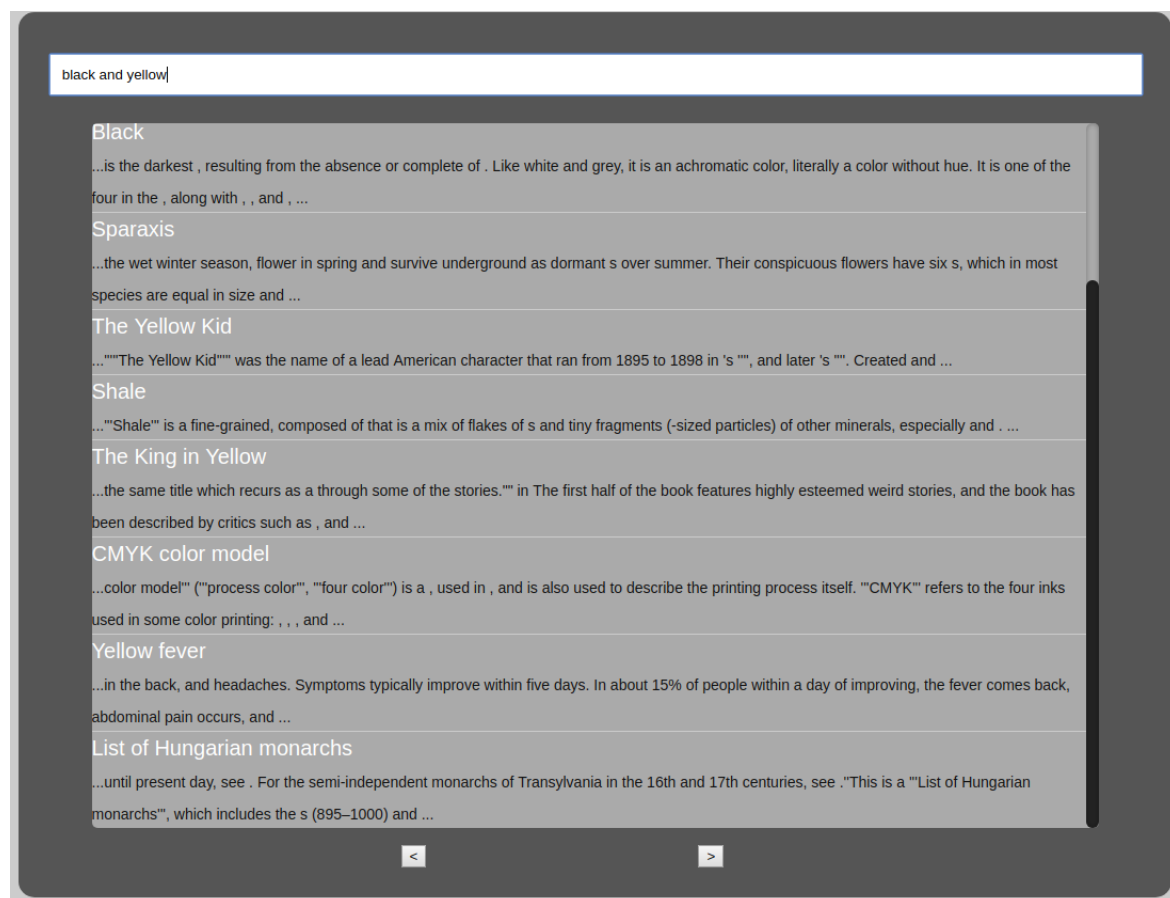
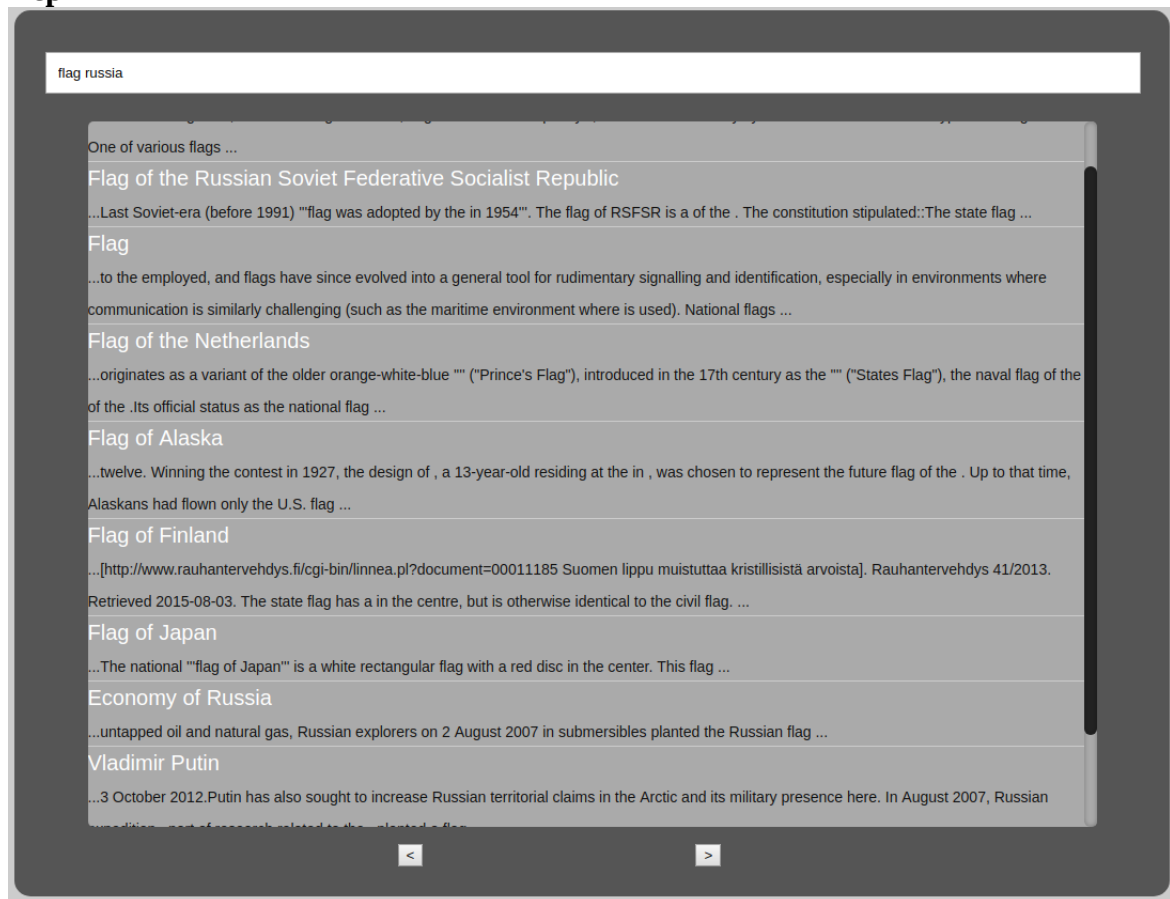
            if mx <= curans:
                mx = len(ans)-1
                ans = spt[:]

            term = ""
        else:
            term += text[i]
            i += 1

    spt = ".join(ans)
    if spt:
        return SNIPPET_FMTSTR.format(spt.strip())
    return "
```

Так же была немного переделана страничка отображения результатов, добавлены сниппеты.

Примеры



Выводы

Построение сниппетов не очень сложная, но крайне полезная задача. Они позволяют довольно точно определять релевантен ли результат поиска и нужно ли переходить к самому документу.