

Подсветка важного в текстах электронной почты

Автор: Кирилл Булгаков, 471 гр.

Научный руководитель: к. т. н., доц. Т.А. Брыксин

Рецензент: ведущий разработчик ООО "Яндекс Технологии"

А.С. Селиванов

Консультант: аналитик ООО "Интеллиджей Лабс" Н.И. Поваров

Санкт-Петербургский Государственный Университет
Кафедра системного программирования, 2019

Проблема

Олимпиада по Теории игр

... ↑ ↓ ×

ки **Олимпиада Теория игр**  olympgamehse@gmail.com
Вам:  thesalnrandom@yandex.ru ^

сегодня в 14:02

← Ответить → Переслать  Удалить  Это спам! ... Ещё

Добрый день, уважаемые студенты!

Спешим информировать о том, что **15-17 марта** состоится **Олимпиада по Теории игр** на базе НИУ ВШЭ, олимпиада пройдет **одновременно в 28 городах** (<https://olymp.hse.ru/ma/2019/place>).

В этом году она пройдет в один тур, **подробнее информацию про олимпиаду можно найти по ссылке** <https://www.hse.ru/news/admission/240281993.html>

Дипломанты олимпиады по этому профилю получают преференции при поступлении на семь магистерских программ, в том числе две — в нижегородской **Вышке** и одну — в **питерской**.

Демо вариант можно найти по ссылке - <https://olymp.hse.ru/data/2018/12/17/1144591184/+390.pdf>

Задания прошлого года можно найти по ссылке

-<https://olympgame.hse.ru/>

Регистрация на Олимпиаду продлится до 4 марта, успевайте!

Если у Вас есть вопросы, обращайтесь!

P.S. передайте, пожалуйста, информацию всем заинтересованным лицам

С уважением,
организационный комитет Олимпиады.

 Написать быстрый ответ всем участникам переписки

Постановка задачи

Цель

- Реализовать решение для подсветки важного в текстах электронной почты на платформе браузерных расширений

Задачи

- Выбрать алгоритмы и определить метрики качества
- Сравнительный анализ алгоритмов
- Спроектировать архитектуру и реализовать решение на платформе браузерных расширений Chromium

Алгоритмы выделения важного

По словам	По предложениям

Алгоритмы выделения важного

По словам	По предложениям
<p data-bbox="467 328 606 366">TF-IDF</p> <ul data-bbox="156 407 942 579" style="list-style-type: none"><li data-bbox="156 407 600 445">• Статистическая мера<li data-bbox="156 476 942 579">• Частота в документе, редкость в корпусе документов	

Алгоритмы выделения важного

По словам	По предложениям
<p data-bbox="467 328 606 366" style="text-align: center;">TF-IDF</p> <ul data-bbox="156 407 942 579" style="list-style-type: none"><li data-bbox="156 407 600 445">● Статистическая мера<li data-bbox="156 476 942 579">● Частота в документе, редкость в корпусе документов	
<p data-bbox="432 721 641 760" style="text-align: center;">Word2Vec</p> <ul data-bbox="156 800 875 907" style="list-style-type: none"><li data-bbox="156 800 627 838">● Векторное разложение<li data-bbox="156 869 875 907">● Норма вектора как степень важности	

Алгоритмы выделения важного

По словам	По предложениям
<p data-bbox="467 328 606 366">TF-IDF</p> <ul data-bbox="156 407 942 579" style="list-style-type: none"><li data-bbox="156 407 600 445">● Статистическая мера<li data-bbox="156 476 942 579">● Частота в документе, редкость в корпусе документов	<p data-bbox="1296 328 1489 366">TextRank</p> <ul data-bbox="1012 407 1715 648" style="list-style-type: none"><li data-bbox="1012 407 1412 445">● Текст в виде графа<li data-bbox="1012 476 1605 514">● Вес ребра - степень важности<li data-bbox="1012 544 1715 583">● Ищем путь с самым большим весом<li data-bbox="1012 613 1514 651">● Готовая реализация: nltk
<p data-bbox="432 738 643 776">Word2Vec</p> <ul data-bbox="156 816 877 921" style="list-style-type: none"><li data-bbox="156 816 629 855">● Векторное разложение<li data-bbox="156 885 877 923">● Норма вектора как степень важности	

Алгоритмы выделения важного

По словам	По предложениям
<p>TF-IDF</p> <ul style="list-style-type: none">• Статистическая мера• Частота в документе, редкость в корпусе документов	<p>TextRank</p> <ul style="list-style-type: none">• Текст в виде графа• Вес ребра - степень важности• Ищем путь с самым большим весом• Готовая реализация: nltk
<p>Word2Vec</p> <ul style="list-style-type: none">• Векторное разложение• Норма вектора как степень важности	<p>LSA</p> <ul style="list-style-type: none">• Текст в виде матрицы• Используем разложение матриц для уменьшения размерности• Готовая реализация: sumy

Оценка качества

- Нет возможности использовать обычные метрики, такие как precision/recall/F-measure и пр.
- Было решено использовать экспертные оценки с использованием сервиса Яндекс.Толока

Платформа Яндекс.Толока

[Машинное обучение, часть 2] Появилось новое домашнее задание

Когда: 25 фев. в 13:50

Кому: mymail@yandex.ru

От кого: noreply@compscicenter.ru

На сайте опубликовано новое задание «Вступительный тест 2 семестр (Машинное обучение, часть 2, весна 2019)», его необходимо сдать до 20:20 11 февраля. Описание домашнего задания ниже, в письме может отображаться некорректно.

Тест по первому семестру

Это письмо отправлено автоматически и не требует ответа.

[Машинное обучение, часть 2] Появилось новое домашнее задание

Когда: 25 фев. в 13:50

Кому: mymail@yandex.ru

От кого: noreply@compscicenter.ru

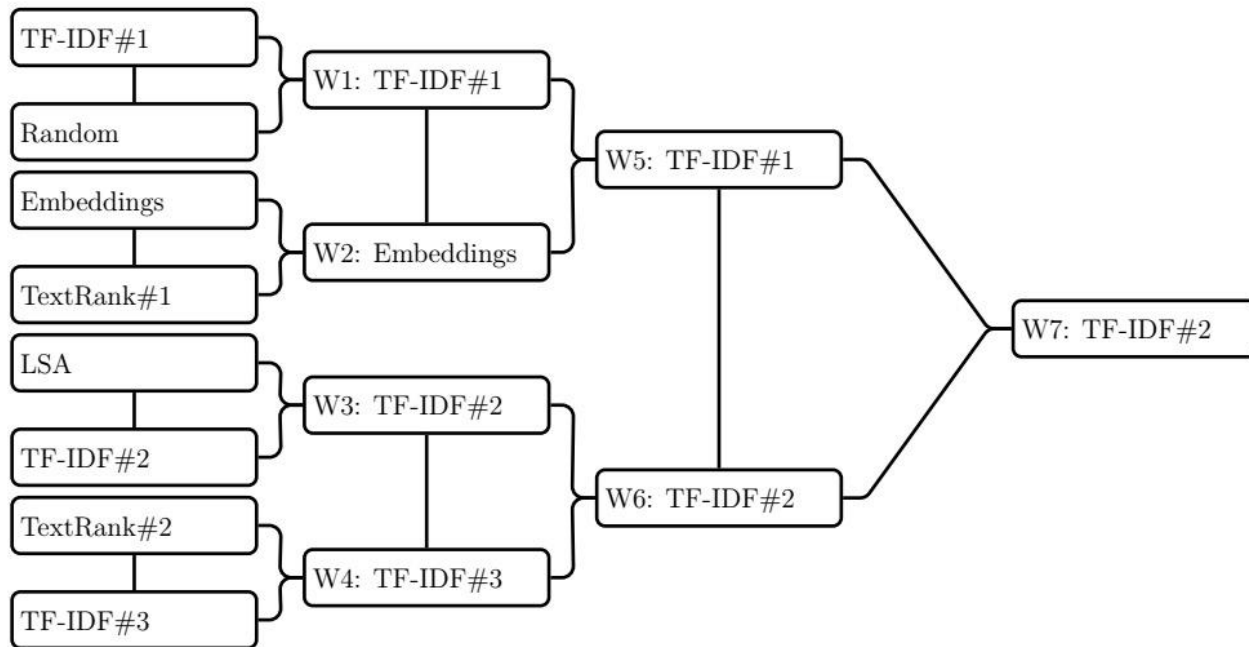
На сайте опубликовано новое задание «Вступительный тест 2 семестр (Машинное обучение, часть 2, весна 2019)», его необходимо сдать до 20:20 11 февраля. Описание домашнего задания ниже, в письме может отображаться некорректно.

Тест по первому семестру

Это письмо отправлено автоматически и не требует ответа.

1 Левый лучше 2 Правый лучше

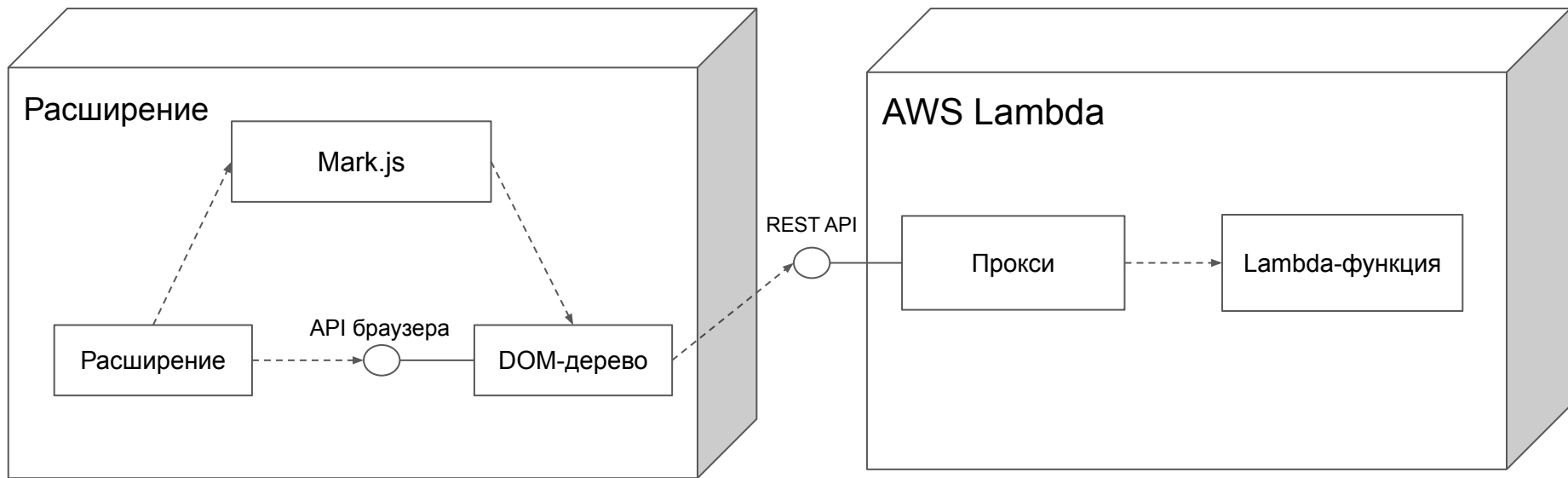
Турнирная сетка для сравнительного анализа алгоритмов



Утилиты для автоматизации сравнительного анализа алгоритмов

- Автоматическое создание скриншотов при помощи Selenium
- Генерация заданий и сбор результатов при помощи API Яндекс.Толока
- Анализ результатов

Архитектура решения



Реализация алгоритмов

- Python
- Упаковка приложения в Docker-контейнер
- Автоматическое развертывание при помощи API AWS Lambda

Результаты

- Выбраны и реализованы следующие алгоритмы для подсветки важного в тексте: TF-IDF в трех вариациях, TextRank в двух вариациях, алгоритм с использованием норм векторов, алгоритм с использованием LSA. Для оценки качества использованы экспертные оценки
- Спроектирован и выполнен на базе платформы Яндекс.Толока эксперимент с целью выявления лучшего алгоритма. Лучшим алгоритмом оказался TF-IDF на национальном корпусе русского языка
- Спроектирована и реализована клиент-серверная система на базе расширений Chromium (использованы технологии AWS Lambda, Docker и языки JavaScript и Python)