

# Curse Words Detector project

Беляев Станислав, Пластинин Виталий

СПб АУ РАН

Весна 2016

# Цели и задачи

- 1 Модуль на питоне для блокировки слов
- 2 Тесты
- 3 Профиллирование работы
- 4 Поместить это внутрь Node.js

# Цели и задачи

- 1 Модуль на питоне для блокировки слов
- 2 Тесты
- 3 Профиллирование работы
- 4 Поместить это внутрь Node.js

# Цели и задачи

- 1 Модуль на питоне для блокировки слов
- 2 Тесты
- 3 Профиллирование работы
- 4 Поместить это внутрь Node.js

# Цели и задачи

- 1 Модуль на питоне для блокировки слов
- 2 Тесты
- 3 Профиллирование работы
- 4 Поместить это внутрь Node.js

# Архитектура проекта

Дано слово  $w$ , нужно найти ближайшее к нему "плохое"  $c$

$$\operatorname{argmax}_c P(c|w) = \operatorname{argmax}_c P(w|c) \frac{P(c)}{P(w)}$$

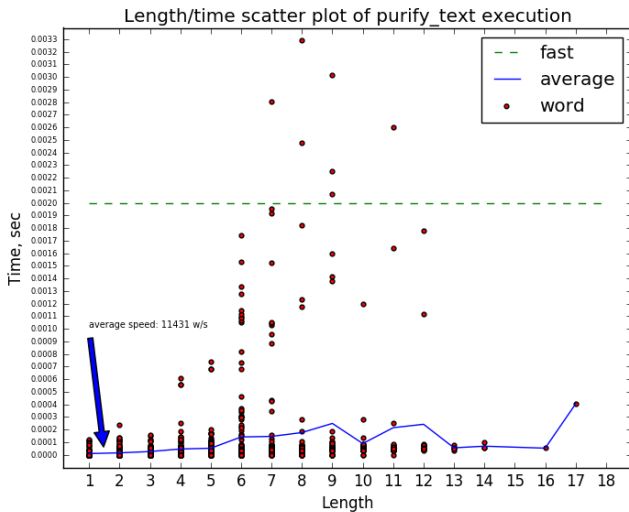
- 1  $P(c)$  - "how likely is  $c$  to appear in text?"
  - частотные списки взял у opencorpora
- 2  $P(w)$  - "same for every  $c$ "
- 3  $P(w|c)$  - вероятность того, что автор написал  $w$ , подразумевая  $c$ 
  - edit distance  $\leq 2$
- 4  $\operatorname{argmax}_c$  - максимум по всем  $c$

# Архитектура проекта

Модули на python3:

- 1 Purifier - main
- 2 Analizier - работа с нормальной формой + предиктор
- 3 Neurister - эвристическая замена цифр и латинских букв на русские
- 4 Statisticer & Tester - классы для сбора статистики и тестов

# Результаты





# Цели и задачи

- 1 База слов
- 2 Рассчёт ближайшего слова
- 3 Замена слов

# Цели и задачи

- 1 База слов
- 2 Рассчёт ближайшего слова
- 3 Замена слов

# Цели и задачи

- 1 База слов
- 2 Рассчёт ближайшего слова
- 3 Замена слов

# Архитектура проекта

- 1 Собираем базу данных слов, так как текущие устарели
  - 1 Обработали все слова на странице
  - 2 Все ссылки положили в очередь, взяли оттуда же следующую страницу
- 2 Дано плохое слово. Подбираем похожее хорошее.

# Архитектура проекта

Модули на python3:

- 1 Sensor - работа со словами
- 2 Crawler - обход сайтов
- 3 DBHelper - для бд

# Примеры

- 1 бл\*дина -> льдина
- 2 ху\*вничать -> чаёвничать
- 3 раз\*банный -> размётанный
- 4 ж\*пы -> окопы

# Конец

Спасибо за внимание

**GitHub:** [github.com/StasBel/CurseWordsDetector](https://github.com/StasBel/CurseWordsDetector)

**GitHub:** [github.com/vitalik239/ClilkBan](https://github.com/vitalik239/ClilkBan)