

# Исправление поисковых запросов

Ильвохин Дмитрий

25 марта 2019 г.

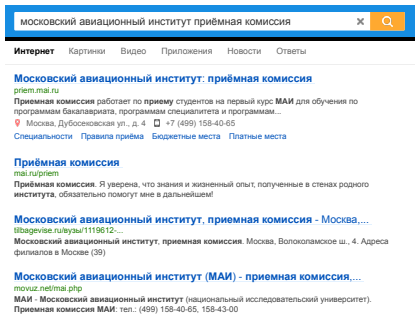
## Запросы с ошибками

Около 11% поисковых запросов содержат ошибки.

Примеры ошибок в одном из самых популярных запросов рунета «вконтакте»:

- ▶ drjynfrnt (154),
- ▶ вкантакте (53),
- ▶ вконтак (34),
- ▶ вконтатке (22),
- ▶ в контакт (18),
- ▶ вконта (14),
- ▶ вконтакты (13),
- ▶ вконтакте (13),
- ▶ в конта (11),
- ▶ вконакте (10).

## Результаты поиска при ошибке в запросе



**Рис. 1:** Результаты по запросу «московский авиационный институт приёмная комиссия»

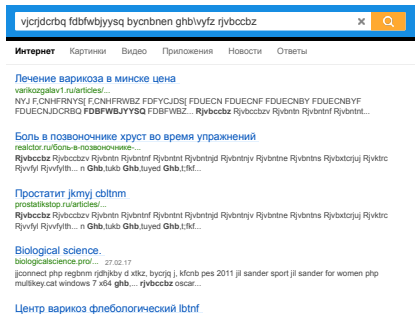


Рис. 2: Результаты по тому же запросу в неверной раскладке клавиатуры

# Типы ошибок

## Замена и пропуск букв

### Замена

- ▶ «вкантакте» → «вконтакте»,
- ▶ «коталог орифлейм» → «каталог орифлейм»,

### Пропуск

- ▶ «однокласники» → одноклассники,
- ▶ «скчать игру» → «скачать игру».

# Типы ошибок

## Вставка и перестановка букв

### Вставка

- ▶ «таныцы айренби видео» → «танцы айренби видео».

### Перестановка

- ▶ «скачтаь медиагет» → «скчатаь медиагет»,
- ▶ «купить барслет» → «купить браслет».

# Типы ошибок

## Разбиение на слова

### Склейка слов

- ▶ «голос 5 сезон бвыпуск» → «голос 5 сезон б выпуск»,
- ▶ «неработаетпробел» → «не работает пробел».

### Расклейка слов

- ▶ «мак бук и вирусы» → «макбук и вирусы»,
- ▶ «ю туб» → «ютуб».

# Типы ошибок

## Раскладка клавиатуры и транслит

### Неверная раскладка клавиатуры

- ▶ «jlyjrkfccybrb» → «одноклассники»,
- ▶ «мл» → «vk».

### Транслитерация

- ▶ «kupit televizor» → «купить телевизор»,
- ▶ «мейл ру» → «mail ru».

# Типы ошибок: сложные случаи

Неправильная раскладка клавиатуры?

vnc → MTC



# Типы ошибок: сложные случаи

Неправильная раскладка клавиатуры?

vnc → MTC

**VNC** — Virtual Network Computing, система удалённого доступа к рабочему столу компьютера.

## Типы ошибок: сложные случаи

Опечатка?

андекс → яндекс

# Типы ошибок: сложные случаи

Опечатка?

андекс → яндекс

**А**ндекс — посёлок в Германии, который известен благодаря находящемуся там бенедиктинскому монастырю.

# Типы ошибок: сложные случаи

Явная ошибка?

пагода → погода

# Типы ошибок: сложные случаи

Явная ошибка?

пагода → погода

**Па́года** — буддийское или индуистское сооружение культового характера.

# Типы ошибок: сложные случаи

Явная ошибка?

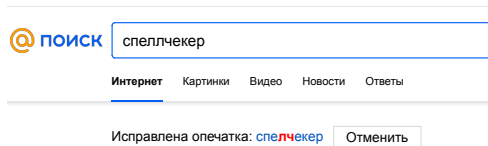


Рис. 3: Большая пагода в Лондоне

# Виды исправлений

## Автоисправление

**Автоисправление** — надежное исправление, в запросе точно допущена ошибка, будут показаны результаты по исправленному запросу.

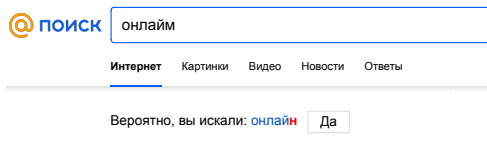


**Рис. 4:** Пример автоисправления: оригинальный запрос: «спелчекер», автоисправление: «спеллчекер»

# Виды исправлений

## Подсказка

**Подсказка** — ненадежное исправление, возможно, в запросе допущена ошибка, будут показаны результаты по оригинальному запросу и предложение исправить запрос.



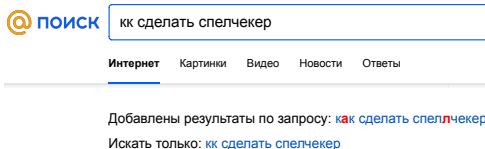
**Рис. 5:** Пример подсказки, оригинальный запрос: «онлайн», подсказка: «онлайн»



# Виды исправлений

## Смешение выдач

**Смешение** — скорее всего в запросе допущена ошибка, будут показаны результаты как по оригинальному запросу, так и по исправленному.



**Рис. 6:** Пример смешения, оригинальный запрос: «кк сделать спелчекер», исправление: «как сделать спеллчекер»

# Модель зашумленного канала<sup>1</sup>

## Задача

- ▶  $\Sigma$  — алфавит;
- ▶  $\Sigma^*$  — множество конечных строк над  $\Sigma$ ;
- ▶  $D \subseteq \Sigma^*$  — словарь корректных слов;
- ▶  $q \notin D$ ,  $c \in D$  — строки над  $\Sigma$ .

## Найти

$$\arg \max_{c \in D} P(c|q).$$

---

<sup>1</sup>Brill E., Moore R. C. An improved error model for noisy channel spelling correction, 2000.

# Модель зашумленного канала

$$\arg \max_{c \in D} P(c|q) = \arg \max_{c \in D} \frac{P(q|c)P(c)}{P(q)} = \arg \max_{c \in D} P(q|c)P(c).$$

## Интерпретация

- ▶  $P(c)$  — модель источника (модель языка),
- ▶  $P(q|c)$  — модель канала (модель ошибок).

## Спеллчекер Питера Норвига (Peter Norvig)<sup>2</sup>

- ▶ Словарь: слова из книг проекта Gutenberg, самые частотные слова из Wiktionary и British National Corpus.
- ▶ Модель языка на основе частот слов в корпусе.
- ▶ Модель ошибок на основе редакционного расстояния.

---

<sup>2</sup><https://norvig.com/spell-correct.html>, 2007.

# Спеллчекер Питера Норвига

## Основная идея

```
def P(word, N=sum(WORDS.values())):  
    return WORDS[word] / N  
  
def correction(word):  
    return max(candidates(word), key=P)  
  
def candidates(word):  
    return (known([word]) or  
            known(edits1(word)) or  
            known(edits2(word)) or [word])  
  
def known(words):  
    return set(w for w in words if w in WORDS)
```

Рис. 7: «Сердце» спеллчекера

## Общая схема работы

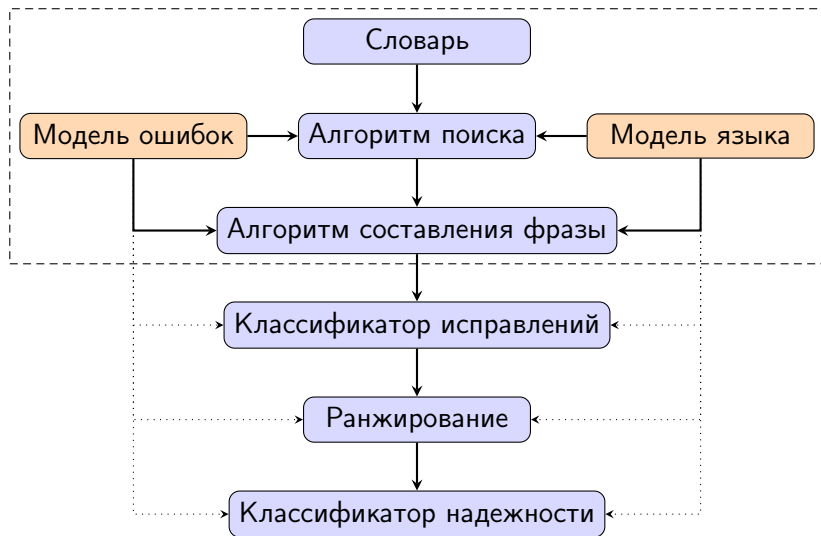


Рис. 8: Составные части системы и их связь

# Модель ошибок<sup>3</sup>

## Идея

$$P(q|c) \propto \text{EditDistance}(q, c)$$

- ▶ Расстояние Левенштейна.
- ▶ Расстояние Дамерау–Левенштейна (транспозиция — одна операция).
- ▶ Взвесить операции редактирования.
- ▶ Вес зависит от контекста и самой операции («тся» → «ться»).

---

<sup>3</sup>Ahmad F., Kondrak G. Learning a spelling error model from search query logs, 2005.

# Модель ошибок

## Сбор данных

Нужны пары *запрос* → *исправление*.

## Использование ассессоров

- ▶ Выдать запрос, попросить исправить.
- ▶ Выдать запрос и исправление, попросить отобрать хорошие.
- ▶ Выдать корректный тест, попросить перепечатать.
- ▶ Выдать аудиофайл, попросить набрать.



# Модель ошибок

## Автоматический сбор данных

Нужны пары *запрос* → *исправление*.

- ▶ Кликовые данные.
  - ▶ Отказы и согласия с исправлениями на страницах с результатами.
  - ▶ Интерливинг.
- ▶ Искусственная генерация данных.

# Модель ошибок

## Построение

Редакционное расстояние между двумя строками  $a$  и  $b$ .

$$\text{lev}_{a,b}(i,j) = \begin{cases} \max(i,j) & \text{if } \min(i,j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1,j) + 1 \\ \text{lev}_{a,b}(i,j-1) + 1 \\ \text{lev}_{a,b}(i-1,j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

- ▶ Используем алгоритм Вагнера–Фишера (Wagner–Fischer).
- ▶ По матрице  $\text{lev}$  получаем редакционное предписание.
- ▶ По редакционному предписанию для всех пар сохраняем статистику.

# Модель ошибок

Построение, модели разных уровней

- ▶ Нулевого уровень: вероятность операции данного типа (вставка, замена, удаление)?
- ▶ Первый уровень: нулевой уровень и дополнительно учитывается текущая буква.
- ▶ Второй уровень: первый уровень и дополнительно учитывается предыдущая буква.
- ▶ ...

# Модель ошибок

## Построение, пример

«президент» → «перзидеед»

Таблица 1: Построение модели ошибок

М	М	D	М	I	М	М	М	М	R	R	М
^	п	р	е	□	з	и	д	е	н	т	\$
^	п	□	е	р	з	и	д	е	е	д	\$

▶ «^п» → «^п» — 1,

▶ ...

▶ «пр» → «п□» — 1,

▶ «т\$» → «д\$» — 1.

# Модель ошибок

Смешение моделей разного уровня

$$\alpha Z + \beta F + \gamma S,$$

где  $Z, F, S$  — модели нулевого, первого и второго уровней,  
 $\alpha, \beta, \gamma \in [0, 1], \alpha + \beta + \gamma = 1$ .

# Модель языка

Пусть  $Q = w_1, w_2, \dots, w_n$  — запрос, состоящий из  $n$  слов,

$$P(Q) = P(w_1)P(w_2|w_1)P(w_3|w_1 w_2) \dots P(w_n|w_1 \dots w_{n-1}).$$

## Униграмная

$$P(Q) \approx P(w_1)P(w_2)P(w_3) \dots P(w_n).$$

## Биграммная

$$P(Q) \approx P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1}).$$

## Техники сглаживания<sup>4</sup>

### Проблема

Пусть  $Q = w_1, w_2, \dots, w_n$  — запрос, состоящий из  $n$  слов, если  $\exists w_k \notin D \Rightarrow P(Q) = 0$ .

---

<sup>4</sup>Chen S. F., Goodman J. An empirical study of smoothing techniques for language modeling, 1999.

# Словарь

## Сбор данных

### Источники данных

- ▶ Словари, книги, тексты новостей.
- ▶ Логи запросов.
- ▶ Слова из веб-документов.



# Алгоритм поиска

## Префиксное дерево

### Идея

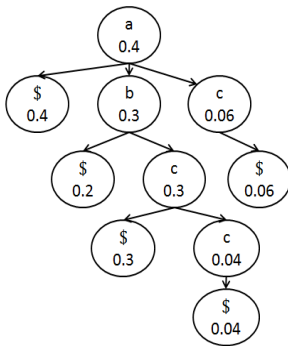
- ▶ Построим префиксное дерево (бор, trie) для слов из словаря.
- ▶ Запустим на нем  $A^*$ .
- ▶ Придумаем эвристику, чтобы быстро работало.

# Префиксное дерево

Эвристики<sup>5</sup>

Query	Count
<i>a</i>	20
<i>ab</i>	10
<i>ac</i>	3
<i>abc</i>	15
<i>abcc</i>	2

a



b

Рис. 9: (a) — набор строк с частотами, (b) — префиксное дерево над этим набором

<sup>5</sup>Duan H., Hsu B. J. P. Online spelling correction for query completion, 2011.



## Идеи для ускорения

- ▶ Используем только  $K$ -лучших переходов для каждой позиции в запросе.
- ▶ Делаем переходы в узлы, если вероятность в них выше определенного процента максимальной вероятности в позиции.
- ▶ Ограничение количества операций редактирования — функция от длины запроса.
- ▶ Предподсчет ответов для частых запросов.

# Алгоритм поиска

Альтернативный подход, фонетические алгоритмы

**Metaphone** — фонетический алгоритм для индексирования слов по их звучанию с учётом основных правил английского произношения.

# Фонетические алгоритмы

Вариант для русского языка<sup>7</sup>

Таблица 2: Аналог Metaphone для русского языка

Код	Буквы	Код	Буквы
1	а, о, ы, у, я	8	г, к, х
3	и, е, ё, ю, я, э	9	л
5	б, п	10	р
6	в, ф	11	м
7	д, т	12	н

Код	Буквы
13	з, с
14	й
15	щ, ч
16	ж, ш
17	ц

<sup>7</sup>Sorokin A. A., Shavrina T. O. Automatic spelling correction for Russian social media texts, 2016.

# Фонетические алгоритмы

## Пример работы

Замены по таблице и несколько других несложных правил.

### Пример

- ▶ «скачать»  $\rightarrow \{13, 8, 1, 15, 1, 7\}$ ,
- ▶ «скочать»  $\rightarrow \{13, 8, 1, 15, 1, 7\}$ .

# Алгоритм составления фразы

## Идея

Пусть исходный запрос  $Q$  был разбит на слова  $w_1, w_2, \dots, w_n$ .  
Для каждого слова были найдены вероятные кандидаты:

- ▶  $w_1 : c_{1,1}, c_{1,2}, \dots, c_{1,k},$
- ▶  $w_2 : c_{2,1}, c_{2,2}, \dots, c_{2,k},$
- ▶  $\dots$
- ▶  $w_n : c_{n,1}, c_{n,2}, \dots, c_{n,k}.$

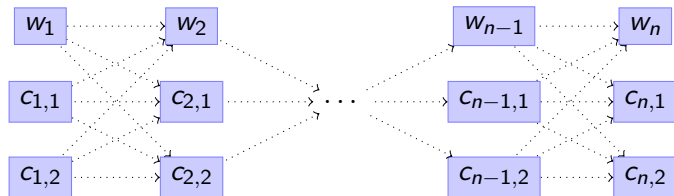


Рис. 11: Схема формирования фразы



# Алгоритм составления фразы

## Имеем

- ▶ Скрытую Марковскую модель и последовательность наблюдений.
- ▶ Веса переходов — вероятности языковой модели  $P(w_i|w_{i-1})$ .

## Найти

Наиболее вероятную последовательность скрытых состояний.

# Алгоритм составления фразы

## Алгоритм Витерби

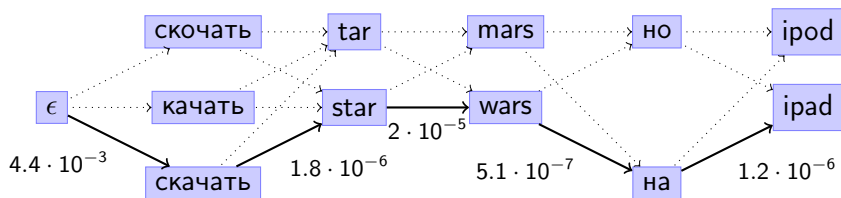


Рис. 12: Пример работы алгоритма

## Биграммная языковая модель

$$viterbi_{i,j} = \max_{0 \leq k < |C_{i-1}|} viterbi_{i-1,k} \times P(c_{i,j} | c_{i-1,k}).$$

## Другие типы исправлений

Общая схема работы похожа на исправление слов.

### Разбиение слова

- ▶ Операции: добавить, убрать разделитель.
- ▶ Оцениваем результат с помощью модели языка.

### Неверная раскладка

- ▶ Собственное разбиение на слова («rf<sub>l</sub>fxjr» → |ка|б|ачок|).
- ▶ Дополнительная операция редактирования: смена раскладки (зависит от типа клавиатуры).

# Классификатор исправлений

- ▶ «нечеткий пиоск строки» → «нечеткий поиск строки»,
- ▶ «нечеткий пиоск строки» ↗ «нечеткий писк строки»,
- ▶ «нечеткий пиоск строки» ↗ «нечеткий киоск строки».

## Признаки

- ▶ Пословные.
- ▶ Позапросные.
- ▶ По типам исправлений (исправление в имени, изменение формы слова).

# Ранжирование

## Сложности

- ▶ Нужно уметь ранжировать разнородные исправления.
- ▶ Сбор обучающих данных.

# Итеративные исправления

Более 90% ошибок исправляются за одну итерацию.

1. «методы государственное подерки литераткра»,
2. «методы государственное подерки литературы»,
3. «методы государственной поддержки литературы».

# Надежность исправлений

Решает нужно ли автоматическое исправление или подсказка.

Точно нужно исправлять, но во что?

- ▶ «поск» → {«поиск», «писк»},
- ▶ «атташа» → {«наташа», «атташе»}.

Вероятность исправления большая, но не факт, что есть ошибка.

- ▶ «vns» → «мтс»,
- ▶ «пагода» → «погода».

## Домашнее задание

Реализовать систему для исправления **слов**.

Kaggle-соревнование

<https://www.kaggle.com/c/itmo-spelling-correction-2019>

Метрика

Среднее расстояние Левенштейна.



# Домашнее задание

## Ограничения

- ▶ Максимальное количество баллов — 25.
- ▶ Решения хуже «No Fixes At All» получают 0 баллов.
- ▶ Не более шести посылок в сутки.
- ▶ Чужие решения сдавать нельзя.
- ▶ Deadline: 04/15/2019 11:59 PM UTC.

# Домашнее задание

## Баллы

Код решения, его короткое описание и kaggle-хэндл нужно прислать на электронную почту.

## Распределение баллов

- ▶ Префиксное дерево и эвристики с лекции — 10 баллов.
- ▶ Фонетические алгоритмы — 3 балла.
- ▶ Модель языка — 2 балла.
- ▶ Модель ошибок — 5 баллов (все уровни).
- ▶ Клёвые эвристики — 5 баллов (не менее двух).

# Домашнее задание

## Контакты

- ▶ Электронная почта: **d.ilvokhin@corp.mail.ru**
- ▶ Telegram: **<https://t.me/r3tsky>**