

Анализ текстов

Лекция

Предобработка и выделение признаков.

Мурат Апишев (great-mel@yandex.ru)

3 октября, 2018

## Предобработка текста

- ▶ Первый шаг любой аналитики – получение данных. Предположим, что данные есть в некотором подходящем для работы формате.
- ▶ Следующая задача – предобработка
- ▶ Базовые шаги предобработки:
  1. токенизация
  2. приведение к нижнему регистру
  3. удаление стоп-слов
  4. удаление пунктуации
  5. фильтрация по частоте/длине/соответствию регулярному выражению
  6. лемматизация или стемминг
- ▶ Чаще всего применяются все эти шаги, но в разных задачах какие-то могут опускаться, поскольку приводят к потере информации.

## Полезные модули

1. `nltk` — один из основных модулей Python для анализа текстов, содержит множество инструментов.
2. `re/regex` — модули для работы с регулярными выражениями
3. `py morphology2/py morphology3` — лемматизаторы
4. Специализированные модули для обучения моделей (например, CRF)
5. `numpy/pandas/scipy/sklearn` — модули общего назначения
6. `codecs` — полезный модуль для работы с кодировками при использовании Python 2.\*

## Токенизация, удаление стоп-слов и пунктуации

В nltk есть разные токенизаторы:

- ▶ RegexpTokenizer
- ▶ BlanklineTokenizer
- ▶ И ещё около десятка штук

Стоп-слова тоже можно удалять с помощью nltk (но лучше дополнительно фильтровать вручную):

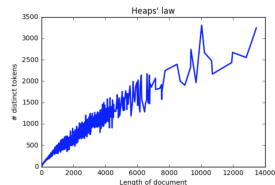
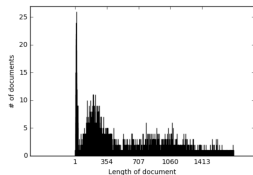
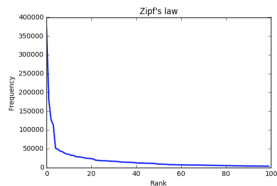
```
1 from nltk.corpus import stopwords
2 words = [word for word in word_list
3           if word not in stopwords.words('russian')]
```

Пунктуацию можно удалять с помощью регулярных выражений, а можно просто:

```
1 from string import punctuation
2 s = ''.join(c for c in s if c not in punctuation)
```

# Статистические свойства коллекции

- ▶ Перед тем, как анализировать коллекцию текстов, необходимо узнать её статистические свойства.
- ▶ Каждая задача требует подсчёта специфичных ей величин.
- ▶ Есть характеристики, на которые нужно смотреть всегда:
  - ▶ Частоты слов в коллекции, закон Ципфа
  - ▶ Гистограмма длин документов
  - ▶ Закон Хипса для документов



# Регулярные выражения

- ▶ Регулярные выражения появились от т.н. регулярных автоматов (классификация грамматик по Хомскому).
- ▶ По факту это некоторый строковый шаблон, на соответствие которому можно проверить текст.
- ▶ Для работы с регулярными выражениями есть множество инструментов и онлайн-сервисов, в Python есть два основных модуля: `re` и `regex`.
- ▶ С синтаксисом можно ознакомиться на странице выбранного инструмента, но основные правила одинаковы, например:
  - ▶ `.` – означает наличие одного любого символа
  - ▶ `[a-zA-Z0-9]` – означает множество символов из заданного диапазона
  - ▶ `+`, `*` – показывают, что следующий перед ними символ или последовательность символов должны повториться  $\geq 1$  раз (`r+`) или  $> 0$  раз (`(xa-)*`)

## Пример из жизни

Регулярные выражения, служащие для определения заголовков публикаций, быстро теряющих актуальность:

```
. *онлайн-| -трансляция. *
```

```
. *в эт(у?и?) минут. *
```

```
. *в эт((от)?и?) час. *
```

```
. *в этот день. *
```

```
. *[0-9]{2}[0-9]{2}( .*|:.*|:. *|;. *|)
```

```
. *[0-9]{1,2}[0-9]{2}([0-9]{2}|[0-9]{4})( .*|:.*|:. *|;. *|)
```

```
. *(от|за|на) [0-9]{1,2} (январ|феврал|март|апрел  
|ию(н|л)|август|сентябр|ноябр|октябр|декабр|ма(я|й)). *
```

```
. *(от|за|на) [0-9]{1,2}[0-9]{2}( .*|:.*|:. *|;. *|)
```

# Стэмминг и лемматизация

**Стэмминг** — процесс приведения слова к основе (отрезание окончания и формообразующего суффикса), грубо, но быстро.

- ▶ Porter stemmer
- ▶ Snowball stemmer
- ▶ Lancaster stemmer

**Лемматизация** — процесс приведения слова к нормальной форме, качественно, но долго.

- ▶ rymorphy2 (язык русский, украинский)
- ▶ mystem3 (язык русский, английский?)
- ▶ Wordnet Lemmatizer (NLTK, язык английский, требует POS метку)
- ▶ Metaphraz (язык русский)
- ▶ Coda/Cadenza (языки русский и английский)



## Пример лемматизации

```
1 import pymorphy2
2
3 text_ru = u'Где твоя ложка, папа?'
4 pymorph = pymorphy2.MorphAnalyzer()
5
6 for word in text_ru.split(u' '):
7     if re.match(u'([a-za-яё]+)', word):
8         word = pymorph.parse(word)[0].normal_form
9     print word,
```

### Вывод:

где твой ложка , папа ?

# Коллокации

**N-граммы** — устойчивые последовательности из N слов, идущих подряд («**машина опорных векторов**»)

**Коллокация** — устойчивое сочетание слов, не обязательно идущих подряд («Он **сломал** своему противнику **руку**»)

**Примеры коллокаций:**

1. *Соединённые Штаты Америки, Европейский Союз*
2. *Машина опорных векторов, испытание Бернулли*
3. *Крепкий чай, крутой кипятилок, свободная пресса*

Часто коллокациями бывают именованные сущности (но далеко не всегда).

## Как можно получать коллокации

- ▶ Извлечение биграмм на основе частот и морфологических шаблонов.
- ▶ Поиск разрывных коллокаций.
- ▶ Извлечение биграмм на основе мер ассоциации и статистических критериев.
- ▶ Алгоритм TextRank для извлечения словосочетаний.
- ▶ Rapid Automatic Keyword Extraction.
- ▶ Выделение ключевых слов по tf-idf.

## Экспериментальные данные

- ▶ Датасет представляет собой статьи о 28 резонансных событиях 2017 года.
- ▶ Каждое событие представлено 100 сырыми текстами.
- ▶ Примеры событий:
  - ▶ *Власти Петербурга согласились передать РПЦ Исаакиевский собор.*
  - ▶ *Дональд Трамп вступил в должность президента США.*
  - ▶ *Умер Дэвид Рокфеллер.*
- ▶ **Ссылка на данные.**
- ▶ Начнём с поиска биграмм в теме «Дональд Трамп вступил в должность президента США»
- ▶ **Ссылка на исходный код** (частично будет ниже).

## Частотные униграммы без стоп-слов

```
1 import nltk
2 import re
3 import pymorphy2
4 from nltk.corpus import stopwords
5
6 prog = re.compile('[А-Яа-я]+')
7 t1 = prog.findall(s.lower())
8
9 morph = pymorphy2.MorphAnalyzer()
10
11 t2 = [morph.parse(token)[0].normal_form
12       for tok in t1
13       if not tok in stopwords.words('russian')]
14 t3 = nltk.FreqDist(t2)
15 t3.most_common(20)
```

## Результат

('трамп', 595)

('президент', 491)

('год', 441)

('который', 428)

('инаугурация', 358)

('сша', 352)

('дональд', 316)

('один', 284)

('россия', 251)

('наш', 223)

('январь', 212)

('это', 198)

('российский', 192)

('время', 184)

('свой', 179)

('быть', 179)

('страна', 173)

('статья', 161)

('человек', 140)

('день', 133)

# Частотные биграммы

- ▶ Без лемматизации и удаления стоп-слов
- ▶ Без POS-тегов, мер ассоциации и т.п.

```
1 bg = list(nltk.bigrams(prog.findall(s.lower())))
2 bgfd = nltk.FreqDist(bg)
3 bgfd.most_common(18)
```

## Результат:

(( 'дональд', 'трамп'), 165)	(( 'по', 'делу'), 50)
(( 'дональда', 'трампа'), 133)	(( 'в', 'должность'), 47)
(( 'президента', 'сша'), 125)	(( 'об', 'этом'), 46)
(( 'в', 'году'), 87)	(( 'в', 'вашингтоне'), 45)
(( 'в', 'россии'), 68)	(( 'из', 'за'), 45)
(( 'избранного', 'президента'), 59)	(( 'в', 'отношении'), 45)
(( 'го', 'президента'), 55)	(( 'президент', 'сша'), 44)
(( 'инаугурация', 'трампа'), 55)	(( 'и', 'в'), 40)
(( 'москва', 'января'), 51)	(( 'млрд', 'руб'), 40)

# Частотные биграммы

- ▶ Без лемматизации и удаления стоп-слов
- ▶ С морфологическим шаблоном (например, Томита)

S -> Adj<gnc-agr[1]> Noun<gnc-agr[1], rt>;

```
1 # s - list with collocation find by Tomita
2 d1 = nltk.FreqDist(s)
3 d1.most_common(16)
```

## Результат:

('избранный президент', 87)  
( '-ый президент', 70)  
( 'белый дом', 69)  
( 'прямая трансляция', 43)  
( 'наша страна', 31)  
( 'этот год', 31)  
( 'соединенный штат', 28)  
( 'новый президент', 27)

('конституционный суд', 25)  
( 'прошлый год', 23)  
( 'весь мир', 21)  
( 'предвыборная кампания', 21)  
( 'ближайшее время', 21)  
( 'новая администрация', 20)  
( 'опасное вождение', 20)  
( 'демократическая партия', 19)



# Поиск разрывных коллокаций

- ▶ Часто устойчивые словосочетания находятся не рядом.
- ▶ **Примеры:**
  - ▶ She *knocked* on his *door*.
  - ▶ They *knocked* on his heavy *door*.
  - ▶ A man *knocked* on the metal front *door*.
- ▶ **Что делаем:**
  - ▶ Рассмотрим все пары слов в некотором окне.
  - ▶ Посчитаем расстояние между словами.
- ▶ **Что меряем:**
  - ▶ **Матожидание** – показывает, насколько часто слова встречаются вместе.
  - ▶ **Дисперсия** – вариабельность позиции.
- ▶ Важно провести лемматизацию.
- ▶ Презентация по теме: <http://tpc.at.ispras.ru/wp-content/uploads/2011/10/lecture4-2016.pdf>

## Пример

$s$	$\bar{d}$	Count	Word 1	Word 2
0.43	0.97	11657	New	York
0.48	1.83	24	previous	games
0.15	2.98	46	minus	points
0.49	3.87	131	hundreds	dollars
4.03	0.44	36	editorial	Atlanta
4.03	0.00	78	ring	New
3.96	0.19	119	point	hundredth
3.96	0.29	106	subscribers	by
1.07	1.45	80	strong	support
1.13	2.57	7	powerful	organizations
1.01	2.00	112	Richard	Nixon
1.05	0.00	10	Garrison	said

Большое значение дисперсии говорит о том, что словосочетание не слишком интересно.

$$\bar{d} = \frac{\sum_{i=1}^n d_i}{n}$$

$$s^2 = \frac{\sum_{i=1}^n (d_i - \bar{d})^2}{n - 1}$$

- ▶  $n$  – число раз, когда два слова встретились.
- ▶  $d_i$  – смещение между словами (может быть  $< 0$ ).
- ▶  $\bar{d}$  – выборочное среднее смещений.

Источник примера

## Меры ассоциации биграмм – PMI

*PMI (Pointwise Mutual Information):*

$$\text{PMI}(w_1, w_2) = \log \frac{f(w_1, w_2)}{f(w_1)f(w_2)}$$

где  $w_i$  – слово,  $f(\cdot)$  – частота слова или биграммы.

- ▶ Оценивает независимость совместного появления пары слов.
- ▶ Значения величины зависят от размеров корпуса.
- ▶ Завышает значимость редких словосочетаний.  
**Решение:** порог по частоте.
- ▶ Выделяет терминологические словосочетания.

## Меры ассоциации биграмм – T-Score

$$\text{T-Score}(w_1, w_2) = \frac{f(w_1, w_2) - f(w_1)f(w_2)}{\sqrt{f(w_1, w_2)/N}}$$

где  $N$  – общее количество биграмм.

- ▶ Является модифицированным ранжированием по частоте.
- ▶ Не преувеличивает значимость редких коллокаций ( $\Rightarrow$  нет необходимости в пороге).
- ▶ Выделяет общеязыковые устойчивые сочетания.
- ▶ По-сути – статистический тест Стьюдента, проверяется гипотеза независимой встречаемости двух слов.

## Меры ассоциации биграмм – T-Score

$t$	$C(w^1)$	$C(w^2)$	$C(w^1 w^2)$	$w^1$	$w^2$
4.4721	42	20	20	Ayatollah	Ruhollah
4.4721	41	27	20	Bette	Midler
4.4720	30	117	20	Agatha	Christie
4.4720	77	59	20	videocassette	recorder
4.4720	24	320	20	unsalted	butter
2.3714	14907	9017	20	first	made
2.2446	13484	10570	20	over	many
1.3685	14734	13478	20	into	them
1.2176	14093	14776	20	like	people
0.8036	15019	15629	20	time	last

Источник примера

## Меры ассоциации биграмм – $\chi^2$ , LLR

$\chi^2$  и LLR (*Log-Likelyhood Ratio*):

- ▶ Так же представляют собой статистические тесты.
- ▶  $\chi^2$  сравнивает наблюдаемые частоты в корпусе с ожидаемыми при верной гипотезе о независимости, большое различие приводит к опровержению гипотезы.
- ▶ LLR: насколько более правдоподобна одна гипотеза, чем другая:
- ▶  $\chi^2$  требует большую выборку наблюдений.
- ▶ **Здесь** можно найти подробные формулы и описания.

## Пример использования

Все описанные меры реализованы в `nltk.collocations`:

- ▶ `bigram_measures.pmi`
- ▶ `bigram_measures.student_t`
- ▶ `bigram_measures.chi_sq`
- ▶ `bigram_measures.likelihood_ratio`

Код, данные и результаты экспериментов доступны [здесь](#).

### Предобработка текстов:

1. Объединим все 2800 текстов в один.
2. Приведём всё к нижнему регистру.
3. Лемматизируем.
4. Удалим стоп-слова.

## Пример

```
1 # m - linear list of tokens
2
3 from nltk.collocations import *
4 N_best = 100 # number of bigrams to extract
5
6 # class for association measures
7 bm = nltk.collocations.BigramAssocMeasures()
8
9 # class for bigrams extraction and storing
10 f = BigramCollocationFinder.from_words(m)
11
12 # remove too seldom bigrams
13 f.apply_freq_filter(5)
```



## Пример

```
1 # get top-100 bigrams using simple frequency
2 raw_freq_ranking = [' '.join(i) for i in
3                     f.nbest(bm.raw_freq, N_best)]
4
5 # get top-100 bigrams using described measures
6 tscore_ranking = [' '.join(i) for i in
7                   f.nbest(bm.student_t, N_best)]
8
9 pmi_ranking = [' '.join(i) for i in
10                f.nbest(bm.pmi, N_best)]
11
12 llr_ranking = [' '.join(i) for i in
13                f.nbest(bm.likelihood_ratio, N_best)]
14
15 chi2_ranking = [' '.join(i) for i in
16                 f.nbest(bm.chi_sq, N_best)]
```

# Результаты

	raw_freq	pmi	t-score	chi2	llr
0	vladimir putin	азотосодержимый добавка	vladimir putin	азотосодержимый добавка	vladimir putin
1	прямая линия	аугусто пиночет	прямая линия	алый парус	прямая линия
2	дональд трамп	бактериальный инфекция	дональд трамп	алёсс капут	риа новость
3	курортный сбор	визитный карточка	курортный сбор	аттестат зрелость	дональд трамп
4	риа новость	висок ром	риа новость	аугусто пиночет	курортный сбор
5	премьер министр	воротничок автоматизация	премьер министр	бактериальный инфекция	исаакиевский собор
6	исаакиевский собор	греф предречь	исаакиевский собор	визитный карточка	премьер министр
7	евгений евтушенко	дилер роад	евгений евтушенко	висок ром	санкт петербург
8	алексей навалный	добавление рацион	алексей навалный	воротничок автоматизация	евгений евтушенко
9	президент сша	доза азотосодержимый	президент сша	генри киссинджер	виталий чуркин
10	чемпионат мир	консультант кредитовать	чемпионат мир	греф предречь	денис вороненков

# TextRank

**TextRank** – приложение алгоритма PageRank к задачам NLP:

- ▶ Строим граф на основе исходного текста
  - ▶  $V$  – вершины (слова)
  - ▶  $E$  – рёбра (связи)
- ▶ Вычисляем веса вершин по мерам центральности, например, по *PageRank*:

$$PR(V_i) = (1 - d) + d \sum_{V_j \in In(V_i)} \frac{PR(V_j)}{Out(V_j)}$$

- ▶ Извлекаем цепочки с наибольшими весами.

Источник: <http://koost.eveel.ru/science/CSEDays2012.pdf>

# Описание TextRank

- ▶ В качестве  $V$  можно взять все уникальные леммы текста (допустимо ограничиться прилагательными и существительными: термины в основном являются именными группами).
- ▶ Сканируем текст с окном из  $N \in [2, 10]$  слов.
- ▶ На каждой итерации считаем для пары слов величину связи

$$WC(w_1, w_2) = \begin{cases} 1 - \frac{d(w_1, w_2) - 1}{N - 1}, & \text{if } d(w_1, w_2) \in (0, N), \\ 0, & \text{if } d(w_1, w_2) \geq N, \end{cases}$$

где  $w_i$  – слова,  $d(w_1, w_2)$  – расстояние между ними (можно просто взять модуль разности позиций).

- ▶ Основание подобной связи – между двумя рядом стоящими словами часто существует семантическое отношение.

# Описание TextRank

- ▶ Ранжируем вершины графа на основании значения TextRank, получаемого случайным блужданием для каждой вершины  $t \in V$ :

$$\text{TR}(t_i) = (1 - d) + d \sum_{t_j \in \text{In}(t_i)} \frac{w_{ji}}{\sum_{t_k \in \text{Out}(t_j)} w_{jk}} \text{TR}(t_j)$$

где  $d$  – фактор затухания,  $\text{In}(t)$  – вершины, входящие в  $t$ ,  $\text{Out}(t)$  – выходящие из  $t$ ,  $w_{ij}$  – вес соответствующего ребра.

- ▶ Упорядочиваем вершины по TR и отбираем  $T$  самых лучших (например,  $T = 1/3|V|$ ). Это множество кандидатов  $C$ .
- ▶ Извлекаем из текста все последовательности слов, состоящие из элементов множества  $C$ .

**Важно:**

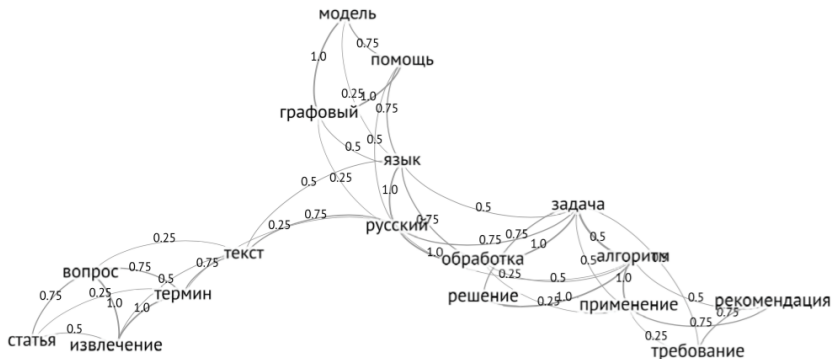
1. в последовательности должно быть хоть одно существительное;
2. в случае вложенности надо рассматривать только последовательность с бльшим весом

# TextRank: пример

## Текст:

Статья посвящена вопросу извлечения терминов из текстов на русском языке при помощи графовых моделей экспериментально исследован алгоритм решения данной задачи. Сформулированы требования и рекомендации к применению алгоритма в задачах обработки русского языка.

## Граф:



## TextRank: пример

### Множество кандидатов:

- ▶ задача – 0,094
- ▶ русский – 0,084
- ▶ алгоритм – 0,083
- ▶ язык – 0,076
- ▶ извлечение – 0,064
- ▶ обработка – 0,063
- ▶ термин – 0,062
- ▶ вопрос – 0,060

### Выделенные словосочетания:

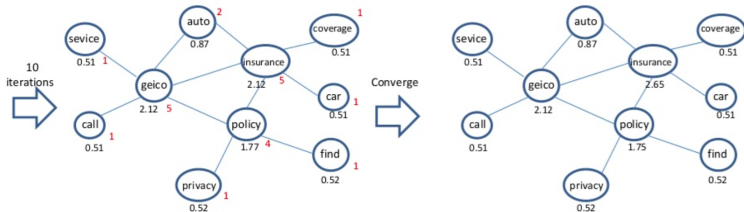
Статья посвящена ***вопросу извлечения терминов*** из текстов на русском языке при помощи графовых моделей экспериментально исследован ***алгоритм*** решения данной ***задачи***. Сформулированы требования и рекомендации к применению ***алгоритма*** в ***задачах обработки русского языка***.

# TextRank: ещё пример

Ссылка на источник

$$S(V_i) = (1 - d) + d * \sum_{j \in \text{nbr}(V_i)} \frac{1}{|\text{degree}(V_j)|} S(V_j)$$

$d$  is the damping factor that usually set to 0.85



**Converge Really Quick!**  
( $\leq 20$  iterations)

$$\text{Score}(\text{geico}) = 0.15 + 0.85 * \left( \underbrace{\frac{1}{1} * 0.51}_{\text{service}} + \underbrace{\frac{1}{1} * 0.51}_{\text{call}} + \underbrace{\frac{1}{2} * 0.87}_{\text{auto}} + \underbrace{\frac{1}{5} * 2.12}_{\text{insurance}} + \underbrace{\frac{1}{4} * 1.77}_{\text{policy}} \right) = 2.12$$

$$\text{Score}(\text{policy}) = 0.15 + 0.85 * \left( \underbrace{\frac{1}{1} * 0.52}_{\text{find}} + \underbrace{\frac{1}{1} * 0.52}_{\text{privacy}} + \underbrace{\frac{1}{5} * 2.12}_{\text{insurance}} + \underbrace{\frac{1}{5} * 2.12}_{\text{geico}} \right) = 1.75$$

$$\text{Score}(\text{service}) = 0.15 + 0.85 * \left( \underbrace{\frac{1}{5} * 2.12}_{\text{geico}} \right) = 0.51$$



# TextRank в Gensim

Полный код и данные доступны [тут](#).

```
1 # text is a string with lemmatized tokens
2 from gensim.summarization import keywords
3 kw = keywords(text)
```

## Результат (примеры):

- ▶ дуров
- ▶ telegram пока
- ▶ роскомнадзор
- ▶ дать компания реестр
- ▶ заявление
- ▶ канал мессенджер
- ▶ возражать против
- ▶ блокировка
- ▶ россия создатель

# RAKE

*RAKE (Rapid Automatic Keyword Extraction):*

- ▶ Фразы-кандидаты – все слова между разделителями.
- ▶ Некоторым образом производится оценка фразы.
- ▶ Фраза-кандидат ограничивается по частоте и количеству слов.
- ▶ Модули: `rake_nltk`, `Rake`, `rake`.
- ▶ [Ссылка](#) на данные и код.

**Пример:**

```
1 from rake_nltk import Rake
2 r = Rake(stopwords.words('russian') + ['это', 'вне'])
3 r.extract_keywords_from_text(tokenized_text)
4 r.get_ranked_phrases()
```

## Выделение ключевых слов по tf-idf

- ▶ **Идея:** хотим выделить слова, которые часто встречаются в данном тексте, и редко – в других текстах.

$$v_{wd} = tf_{wd} \times \log \frac{N}{df_w}$$

где  $tf_{wd}$  – число раз, которое слово  $w$  встретилось в документе  $d$ ,  $df_w$  – число документов, содержащих  $w$ ,  $N$  – общее число документов.

- ▶ Такие слова, как правило, информативны, и значение tf-idf является хорошим признаком.
- ▶ Значения tf-idf для слов текста можно получить с помощью `sklearn.feature_extraction.text.TfidfVectorizer`
- ▶ **Ссылка** на данные и код примера.