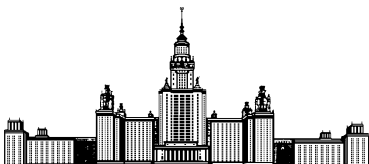


# Ранжирование текстовых документов на основе оценок когнитивной сложности текста

**Выполнил: Еремеев Максим Алексеевич, 417 гр.**

Научный руководитель: д.ф.-м.н. Воронцов Константин Вячеславович  
(ММП - ВМК - МГУ им. М.В. Ломоносова)



Москва • 11 июня 2020

## Индексы удобочитаемости текста

*Индексы удобочитаемости* — меры сложности восприятия текста

- Наиболее распространенные:
  - индекс Флеша,
  - *Automated Readability index (ARI)*,
  - *SMOG-index*
- Вычисляются на основе простых параметров:
  - длины предложений и слов
  - среднего количества букв в слове
  - среднего количества слов в предложении
- содержат большое количество констант, которые подбираются отдельно для каждого языка

$$\text{Flesh}(d) = 206.835 - 1.015 \times \frac{w}{s} - 84.6 \times \frac{\text{syl}}{w}$$

---

*Flesh, R.* How To Test Readability. 1951.

*Senter, R.J. and Smith, E.A.* Automated Readability Index. 1967.

## Ранжирование документов в порядке чтения

*Reading Order* — порядок документов от общих и простых к сложным и узкоспециализированным.

### Применения:

- Сортировка выдачи разведочного информационного поиска
- Персонализация образовательных траекторий

### Подход, основанный на тематическом моделировании:

- Оценка общности текста — энтропия тематического вектора
- Оценка близости документов — косинусное расстояние между тематическими эмбедами
- В цепочку попадают близкие документы в порядке убывания общности

---

*Georgia Koutrika, Lei Liu, Steven Simske. Generating Reading Orders over Document Collections. 2015.*

## Психофизиологическая нагрузка декодирования речи

Что такое «тяжёлый текст» и чем он отличается от лёгкого?

**Пример.** Частота буквы **р** в русском языке 0.04, но здесь 0.17:



**Гипотезы** из нейрофизиологии и психофизиологии речи:

- декодирование каждого элемента языка вызывает нагрузку определенной зоны коры головного мозга
- зоне требуется *время рефрактерности* для восстановления
- в ходе эволюции языка устанавливаются существенно неравномерные распределения частот (закон Ципфа)
- в ходе освоения языка для высокочастотных элементов устанавливаются более короткие периоды рефрактерности

А.А.Биркин. Код речи. СПб.: Гиппократ, 2007.

## Идея комплексного оценивания сложности текста

### Основные положения предлагаемого подхода:

- *уровни языка*: фонетический, морфологический, лексический, синтаксический, дискурсивный
- на уровне  $h$  текст представляется в виде последовательности *токенов* алфавита  $A_h$
- *сложность текста* на уровне  $h$  — это доля токенов, имеющих аномально высокую нагрузку
- нагрузка токена *аномально высокая*, если она превышает 95%-ю квантиль его нагрузки в референтном корпусе
- *референтный корпус* — тексты, которые можно считать простыми для выбранной читательской аудитории

---

*M.Eremeev, K.Vorontsov. Lexical quantile-based text complexity measure. RANLP-2019.*

## Квантильный подход к оцениванию сложности текста

$x_1, \dots, x_n$  — последовательность токенов из  $A_h$  в тексте  $d$ ;

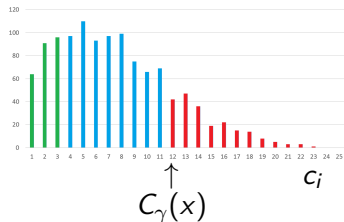
$c_i$  — нагрузка (оценка сложности) токена  $x_i$ ;

$w_i$  — вес нагрузки токена  $x_i$ ;

Оценка сложности текста — суммарный вес сложных токенов:

$$W(d) = \sum_{i=1}^n w_i [c_i > C_\gamma(x_i)]$$

$C_\gamma(x)$  —  $\gamma$ -квантиль распределения сложности токена  $x$  в референтном корпусе несложных текстов



M.Eremeev, K.Vorontsov. Lexical quantile-based text complexity measure. RANLP-2019.

## Два типа оценок токенов — частотные и сложностные

### Частотная оценка нагрузки токена

- для каждого токена  $a \in A_h$  строится эмпирическое распределение значений его нагрузок  $\{c_i : x_i = a\}$  по референтному корпусу несложных текстов
- $c_i = f(r_i)$  — убывающая функция расстояния  $r_i$  от токена  $x_i$  до его предыдущего вхождения
- примеры: частоты букв,  $n$ -грамм, слов

### Сложностная оценка нагрузки токена

- одноэлементный алфавит  $A_h = \{a\}$ , строится эмпирическое распределение всех нагрузок  $\{c_i\}$  по референтному корпусу несложных текстов
- $c_i$  — числовая характеристика сложности элемента  $x_i$
- примеры: длина слова, предложения, синтаксической связи

## Фонетический уровень: токены — это буквы

$r_i$  — расстояние от токена  $x_i$  до его предыдущего вхождения:

$$\dots \boxed{x_{i-r_i} = a} \underbrace{x_{i-r_i+1} \ x_{i-r_i+2} \ \dots \ x_{i-2} \ x_{i-1} \ \boxed{x_i = a}}_{r_i} \dots$$

Если предыдущего вхождения нет, доопределяем  $r_i$  «через хвост» (тогда сумма  $r_i$  по всем вхождениям  $x_i = a$  равна длине текста  $n$ ):

токен	г	е	р	о	й	н	а	ш	е	г	о	в	р	е	м	е	н	и
$r_i$ исходное	–	–	–	–	–	–	–	–	7	9	7	–	10	5	–	2	11	–
$r_i$ доопред.	9	4	8	11	18	7	18	18	7	9	7	18	10	5	18	2	11	18

Пример определения  $c_i$  как убывающей функции от  $r_i$ :

$$c_i = 1/r_i$$



# Морфологический уровень: токены — это буквенные $n$ -граммы

Варианты определения токенов:

- слоги,
- буквенные  $n$ -граммы,
- буквенные  $n$ -граммы, не сохраняющие порядок букв,
- морфемы (приставки, корни, суффиксы, окончания)

**Пример.**

1-граммы	г	е	р	о	й	н	а	ш	е	г	о	в	р	е	м	е	н	и
----------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$n$ -граммы, сохраняющие порядок букв

2-граммы	ге	ер	ро	ой		на	аш	ше	ег	го		вр	ре	ем	ме	ен	ни	
3-граммы	гер	еро	рой			наш	аше	шег	его			вре	рем	еме	мен	ени		

$n$ -граммы, не сохраняющие порядок букв

2-граммы	ге	ер	ор	йо		ан	аш	еш	ге	го		вр	ер	ем	ем	ен	ин	
3-граммы	гер	еор	йор			анш	аше	геш	гео			вер	емр	еем	емн	еин		

## Лексический уровень: токены — это слова

### Варианты определения токенов для частотных оценок:

- слова после лемматизации,
- словные  $n$ -граммы,
- словные  $n$ -граммы, не сохраняющие порядок слов,

### Варианты сложностных оценок нагрузки:

- $c_i$  — длина слова
- $c_i = \frac{1}{\text{count}(x_i)}$  — редкость слова в референтном корпусе,  
где  $\text{count}(x_i)$  — количество вхождений слова

## Синтаксический уровень

Используется синтаксический парсер UDPipe. Токенами являются предложения.

### Варианты определения токенов для частотных оценок:

- грамматические структуры — *синтагмы*, в которых отброшены слова и оставлены теги частей речи и/или членов предложения

### Варианты сложностных оценок нагрузки:

- $c_i$  — максимальная длина синтаксической связи в предложении
- $c_i$  — суммарная длина синтаксических связей
- $c_i$  — средняя длина синтаксических связей

## Дискурсивный уровень

На дискурсивном уровне оценивается осмысленность текста, его связность и последовательность.

### Варианты сложностных оценок нагрузки:

- $c_i$  — число слов в предложении
- $c_i$  — число логических связей в предложении (и, или, значит, который, чтобы, . . . около 150 выражений)

## Обучаемая линейная модель когнитивной сложности текста

Пусть  $W_k(d)$ ,  $k = 1, \dots, K$  — различные оценки сложности.

Линейная агрегированная оценка сложности с параметрами  $\alpha_k$ :

$$W(d, \alpha) = \sum_{k=1}^K \alpha_k W_k(d), \quad \alpha_k \geq 0.$$

Данные экспертного сравнения пар документов:

$d \prec d'$  — документ  $d'$  сложнее документа  $d$ .

Критерий обучения агрегированной оценки:

$$\sum_{d \prec d'} \underbrace{\mathcal{L}(W(d', \alpha) - W(d, \alpha))}_{\text{pair-wise margin}} \rightarrow \min_{\alpha},$$

где  $\mathcal{L}(M)$  — гладкая невозрастающая функция отступа  $M$ .

## Сбор ассессорских оценок

Для оценивания качества была сгенерирована выборка пар статей русскоязычной Википедии из категорий математики, физики, химии, информатики.

Ассессорам предлагалось выбрать из двух статей ту, которая потребовала больше усилий для её понимания и содержала больше незнакомых терминов, либо указать, что статьи примерно равны по сложности, либо что они совершенно из разных областей.

Assiguasy — доля пар, на которых и модель, и ассессоры выбрали одну и ту же статью как сложную.

### Какая из статей сложнее?

The screenshot shows a comparison interface with two article snippets. The left snippet is titled 'Физический вакуум' and discusses concepts like 'абсолютный вакуум' and 'квантовый вакуум'. The right snippet is titled 'Матрица смежности' and discusses graph theory concepts like 'узлы' and 'ребра'. Below the snippets are three buttons: 'Левая', 'Равны', and 'Правая'. At the bottom, there is a button labeled 'Невозможно определить'.

- Левая
- Равны
- Правая
- Невозможно определить

## Эксперимент 1. Отдельные модели

Класс моделей	Модель	Accuracy
Индексы Удобочитаемости	Automated Readability Index	50.5%
	Flesch-Kincaid Grade	44.7%
	Flesch Reading Ease	50.7%
	Coleman-Liau	<b>52.1%</b>
Фонетические	Частотная (буквы)	62.5%
Морфологические	Частотная (с учетом порядка букв)	70.9%
	Частотная (без учета порядка букв)	73.1%
Лексические	Сложностная (длина слова)	42.4%
	Частотная	<b>75.0%</b>
	Сложностная	71.2%
Синтаксические	Сложностная (длина синтаксической связи)	62.0%
	Частотная (синтгамы)	64.2%
Дискурсивные	Сложностная (слова-связки)	62.5%

## Эксперимент 2. Агрегированная модель

Модель	Функция отступа	Accuracy
Coleman-Liau	-	52.1%
Морфологическая (без учета порядка)	-	73.1%
Лексическая (частотная)	-	75.0%
Синтаксическая (синтгамы)	-	64.2%
Дискурсивная	-	62.5%
Агрегированная модель	Отрицательная квадр. ошибка	<b>88.1%</b>
Агрегированная модель	Отрицательная сигмоида	84.6%
Агрегированная модель	Отрицательная абс. ошибка	85.1%

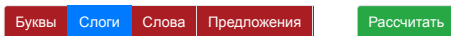
- Агрегирование моделей со всех уровней дает прирост качества
- Все квантильные модели превосходят индексы удобочитаемости
- Частотная лексическая модель показывает лучшее качество среди квантильных моделей



## Веб-ресурс TextComplexity.Net

Интерактивная веб-система подсвечивает сложных токены на каждом уровне в тексте при его модификации с целью упрощения. Ресурс использует оригинальную библиотеку с открытым кодом *Cognitive-Complexity*.

Выберите тип токенов:



Введите текст:

на дворе трава на траве дрова, не руби дрова на траве двора

Результат:

на дво ре тра ва на тра ве дро ва , не руби дро ва на тра ве двора

## Результаты, выносимые на защиту

- *Предложен* квантильный подход к оцениванию сложности текстов, применимый единообразно на всех уровнях языка.
- *Предложена* методика верификации оценок сложности.
- *Предложен* обучаемый метод агрегирования оценок сложности со всех уровней.
- *Разработан* веб-ресурс [TextComplexity.net](http://TextComplexity.net), позволяющий определять сложность произвольного текста.
- *Разработана* библиотека для оценивания сложности текста с открытым кодом [Cognitive-Complexity](https://github.com/TextComplexity/Cognitive-Complexity).
- *Показана* применимость оценок сложности для ранжирования текстовых документов.
- *В экспериментах показана* высокая точность квантильных оценок.