

**Отчёт по практическому заданию  
"Униграммная и биграммная языковые  
модели со сглаживанием"**

Выполнил: Алексей Грищенко, 209 группа

# Постановка задачи

Построить биграммную и униграммную модель и оценить ее качество с помощью перплексии на заданном кусочке текста при различных параметрах сглаживания.

## Используемые формулы

Пусть дан тестовый корпус:  $w_1w_2...w_n$  Вычисление вероятности для униграммы осуществляется по следующей формуле:

$$P(w_i) = \frac{c(w_i) + 1 * \lambda}{N + \lambda * V}$$

где  $c(w_i)$  - частота униграммы в обучающем корпусе,  $N$  - количество слов в обучающем корпусе,  $V$  - количество слов в словаре,  $\lambda$  - параметр сглаживания.

Вычисление вероятности для биграммы вычисляется по следующей формуле:

$$P(w_i|w_{i-1}) = \frac{c(w_iw_{i-1}) + 1 * \lambda}{c(w_{i-1}) + \lambda * V^2}$$

где  $c(w_iw_{i-1})$  - частота подстроки  $w_iw_{i-1}$

Как уже было сказано выше, оценка качества языковой модели будет производиться с помощью перплексии. По сути, этот показатель говорит о том, какое количество вариантов в среднем рассматривает наша языковая модель на каждом шаге.

Перплексия для униграммной и биграммной моделей вычисляется по следующим формулам:

$$P_{uni} = \sqrt[n]{\prod_{i=1}^n 1/P(w_i)}$$
$$P_{bi} = \sqrt[n]{\prod_{i=1}^n 1/P(w_iw_{i-1})}$$

## Алгоритм работы программы

Рассмотрим алгоритм работы биграммной языковой модели. В качестве корпуса берётся стих Самуила Маршака "Дом, который построил Джек". В нём 247 слов, из которых 58 слов являются уникальными. Из этого корпуса выделим два непересекающихся части. Первая часть используется для обучения нашей модели. В ней 37 слов, из которых 21 уникальных и хранится она в файле `train.txt`. Вторая часть хранится в `test.txt` и используется для оценки качества.

Первым делом программа сканирует весь стих, для того, чтобы выделить для себя уникальные слова. Далее программа сканирует обучающий корпус. В процессе сканирования в соответствующих структурах данных происходит заполнение информации о частотах уникальных слов и биграмм, встреченных в обучающем корпусе. На последнем этапе идет вычисление перплексии за счёт полученных данных при различных параметрах лямбда. Аналогичный алгоритм применяется в униграммной языковой модели.

## Результаты

Следующие таблицы отражают влияние сглаживания с параметром  $\lambda = 1$  на подсчёт вероятности некоторых униграмм/биграмм в построенной языковой модели:

Униграмма	Р до сглаживания	Р после сглаживания
Который	0.081	0.042
часто	0.027	0.021
кот	0	0.010
пёс	0	0.010

Биграмма	Р до сглаживания	Р после сглаживания
В доме	1	0.00089
Которая в	0.666	0.00089
Который синица	0	0.00029
В Джек	0	0.00029

Из таблиц видно, что сглаживание с параметром  $\lambda = 1$  является неэффективным, т.к она слишком много вероятности отдаёт униграммам и биграммам, которые никогда не встречались в обучающем корпусе.

Что касается зависимости перплексии от параметра сглаживания, то из нижепредставленных графиков следует, что в униграммной модели наиболее точная модель получается при  $\lambda \approx 0.23$ . В биграммной модели лучший результат получается при близком к нулю параметре лямбда.

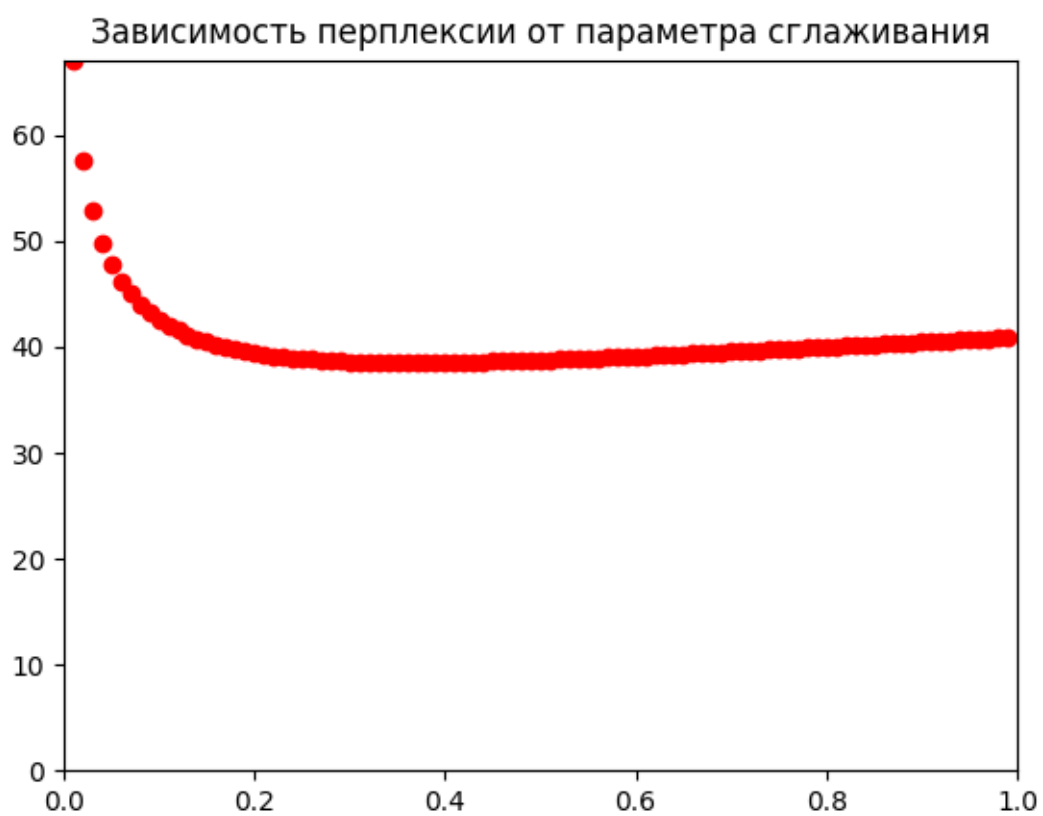


Рис. 1: Зависимость перплексии униграммной модели от параметра сглаживания

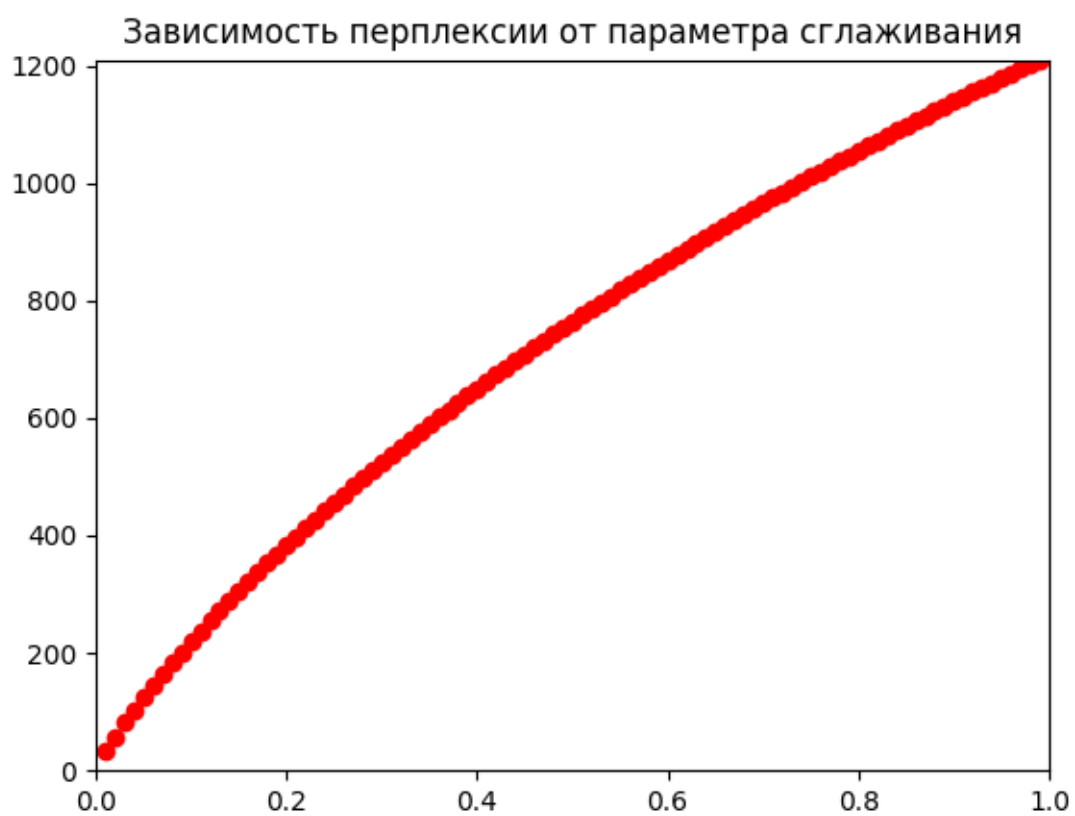


Рис. 2: Зависимость перплексии биграммной модели от параметра сглаживания