



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

Кафедра прикладной математики и
информатики

Компьютерная лингвистика. Практика.

#11, #12, #13, #14, #15, #16

Нижний Новгород, 2018



- ~~Git. Python. Pandas~~
- ~~XML. JSON. Парсинг новостного сайта~~
- ~~Предобработка корпуса. Лемматизация. Дистанция редактирования~~
- ~~Обратный индекс~~
- ~~TF-IDF. Разреженные матрицы~~
- ~~Косинусное расстояние. BIRRM~~
- ~~Анализ моделей классификации. Ошибки 1го, 2го рода, точность, полнота, F-score... (sklearn)~~
- ~~Марковский процесс. N-граммы. Правдоподобие и перплексия. OpenCorpora~~
- **Тегирование и POS. HMM, TreeTagger, Spacy**
- **Тематическое моделирование**
- **Векторное представление слов, skip-gram, w2vec. Введение в нейронные сети**

Практика #15.

Векторное представление слов.

- Skip-Gram

Source Text	Training Samples
<div>The quick brown fox jumps over the lazy dog.</div>	<div>(the, quick)</div> <div>(the, brown)</div>
<div>The quick brown fox jumps over the lazy dog.</div>	<div>(quick, the)</div> <div>(quick, brown)</div> <div>(quick, fox)</div>
<div>The quick brown fox jumps over the lazy dog.</div>	<div>(brown, the)</div> <div>(brown, quick)</div> <div>(brown, fox)</div> <div>(brown, jumps)</div>
<div>The quick brown fox jumps over the lazy dog.</div>	<div>(fox, quick)</div> <div>(fox, brown)</div> <div>(fox, jumps)</div> <div>(fox, over)</div>



- CBOW (Continuous Bag Of Words)

context word center words

I like playing football with my friends →

[0, 1, 0, 0, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0]

[1, 0, 0, 0, 0, 0, 0]

I like playing football with my friends →

[1, 0, 0, 0, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0]
[0, 0, 0, 1, 0, 0, 0]

[0, 1, 0, 0, 0, 0, 0]

I like playing football with my friends →

[1, 0, 0, 0, 0, 0, 0]
[0, 1, 0, 0, 0, 0, 0]
[0, 0, 0, 1, 0, 0, 0]
[0, 0, 0, 0, 1, 0, 0]

[0, 0, 1, 0, 0, 0, 0]

I like playing football with my friends →

[0, 1, 0, 0, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0]
[0, 0, 0, 0, 1, 0, 0]
[0, 0, 0, 0, 0, 1, 0]

[0, 0, 0, 1, 0, 0, 0]

I like playing football with my friends →

[0, 0, 1, 0, 0, 0, 0]
[0, 0, 0, 1, 0, 0, 0]
[0, 0, 0, 0, 0, 1, 0]
[0, 0, 0, 0, 0, 0, 1]

[0, 0, 0, 0, 1, 0, 0]

I like playing football with my friends →

[0, 0, 0, 1, 0, 0, 0]
[0, 0, 0, 0, 1, 0, 0]
[0, 0, 0, 0, 0, 0, 1]

[0, 0, 0, 0, 0, 1, 0]

I like playing football with my friends →

[0, 0, 0, 0, 1, 0, 0]
[0, 0, 0, 0, 0, 1, 0]

[0, 0, 0, 0, 0, 0, 1]



Контрольная работа #4.

Программа классификации и анализа точности классификации

Задача.

1. Используя вектора TF-IDF и LSA, выполнить тематическое моделирование в выбранном корпусе
2. Найти оптимальное количество латентных тем и способ препроцессинга



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ