

Задача 1: сравнение предложений

Дан набор предложений, скопированных с Википедии. Каждое из них имеет "кошачью тему" в одном из трех смыслов:

- кошки (животные)
- UNIX-утилита `cat` для вывода содержимого файлов
- версии операционной системы OS X, названные в честь семейства кошачьих

Ваша задача — найти два предложения, которые ближе всего по смыслу к расположенному в самой первой строке. В качестве меры близости по смыслу мы будем использовать косинусное расстояние.

`sentences.txt`

Выполните следующие шаги:

1. Скачайте файл с предложениями (`sentences.txt`).
2. Каждая строка в файле соответствует одному предложению. Считайте их, приведите каждую к нижнему регистру с помощью строковой функции `lower()`.
3. Произведите токенизацию, то есть разбиение текстов на слова. Для этого можно воспользоваться регулярным выражением, которое считает разделителем любой символ, не являющийся буквой: `re.split('[^a-z]', t)`. Не забудьте удалить пустые слова после разделения.
4. Составьте список всех слов, встречающихся в предложениях. Сопоставьте каждому слову индекс от нуля до $(d - 1)$, где d — число различных слов в предложениях. Для этого удобно воспользоваться структурой `dict`.
5. Создайте матрицу размера $n * d$, где n — число предложений. Заполните ее: элемент с индексом (i, j) в этой матрице должен быть равен количеству вхождений j -го слова в i -е предложение. У вас должна получиться матрица размера $22 * 254$.
6. Найдите косинусное расстояние от предложения в самой первой строке (`In comparison to dogs, cats have not undergone...`) до всех остальных с помощью функции `scipy.spatial.distance.cosine`. Какие номера у двух предложений, ближайших к нему по этому расстоянию (строки нумеруются с нуля)? Эти два числа и будут ответами на задание. Само предложение (`In comparison to dogs, cats have not undergone...`) имеет индекс 0.
7. В качестве ответа предоставьте номера предложений и косинусные расстояния. Совпадают ли ближайшие два предложения по тематике с первым? Совпадают ли тематики у следующих по близости предложений?

Подключение Google-диска к <https://colab.research.google.com/>

- Подключение библиотеки:

```
from google.colab import drive
drive.mount('/content/drive')
```

- Стандартный путь к директории `'drive/My Drive/Colab Notebooks'`

- Пример чтения файла:
`Date= pd.read_csv('drive/My Drive/Colab Notebooks/date1.csv', index_col=False)`

Материалы

Справка по функциям пакета `scipy.linalg`: <http://docs.scipy.org/doc/scipy/reference/linalg.html>

Справка по работе с файлами в Python: <https://docs.python.org/2/tutorial/inputoutput.html#r..>

Справка по регулярным выражениям в Python (если вы захотите узнать про них чуть больше): <https://docs.python.org/2/library/re.html>