

ЛР1: Токенизация

Задание

Реализовать процесс разбиения текстов документов на токены, который потом будет использоваться при индексации. Для этого потребуется выработать правила, по которым текст делится на токены. Необходимо описать их в отчёте, указать достоинства и недостатки выбранного метода. Привести примеры токенов, которые были выделены неудачно, объяснить, как можно было бы поправить правила, чтобы исправить найденные проблемы.

Метод решения

1. Изучение способов токенизации текста с использованием разных библиотек python.
2. Написание и отладка кода, выполняющего разделение текста на отдельные токены.
3. Сбор и анализ статистических данных.

В качестве инструмента по токенизации текста была выбрана библиотека nltk, которая обладает словарем стоп-слов русского языка. Дополнительно текст был очищен от мусорных знаков: знаков пунктуации, знаков операций и других ascii символов, которые не представляют буквы русского алфавита.

Результаты выполнения

Файлы	Размер до токенизации	Размер после токенизации	Средняя длина токена	Время
все статьи в формате json	11Gb	9.9Gb	6.33	1:29:40

Исходные данные

 {{id: "7", "url": "https://ru.wikipedia.org/wiki?curid=7", "title": "Литва", "text": "ЛитваЛитва́ (Литва́), официальное название — Литовская Республика (Литва) — государство, расположенное в северной части Европы. Площадь — км². Протяжённость с севера на юг — 280 км, а с запада на восток — 370 км. Население составляет человек (январь, 2021). Занимает 137-е место в мире по численности населения и 121-е по территории. Имеет выход к Балтийскому морю, расположена на его восточном побережье. Береговая линия составляет всего 99 км (наименьший показатель среди государств Балтии). На севере граничит с Латвией, на юго-востоке — с Белоруссией, на юго-западе — с Польшей и Калининградской областью России. По площади и населению является самым крупным прибалтийским государством.ЛитваСтолица — Вильнюс. Официальный язык — литовский. Де-юре — не признан прибалтийским государством. Независимости страны провозглашено 11 марта 1990 года. 6 сентября 1991 года Государственный совет СССР признал независимость Литвы.ЛитваЛитва — член ООН (1991), ОБСЕ (1991), Совета Европы (1993), ВТО (2001), Европейского союза (2004), НАТО (2004) и ОЭСР (2018). Входит в Шенгенскую зону и Еврозону.ЛитваИтология слова «Литва» точно не известна, при этом существует множество версий, ни одна из которых не получила всеобщего признания. Корень «лит» и его варианты «лет/лёт» допускают различные толкования как в балтских и славянских, так и в других индоевропейских языках. Так, например, существуют созвучные топонимы на территории Словакии («Litvaľ») и Румынии («Lituaľ»), известные с XI—XII веков. По мнению Е. Поспелова, топоним образован от древнего названия реки Летава (Lietavā от «лит», русское «Летаука»). Феодалное княжество, по землям которого протекала эта река, со временем заняло ведущее положение и название было распространено на всё государство. В «Повести временных лет» (XII век) упоминается этноним «литва», полностью совпадающий с названием местности «Литавы» и по смыслу (территория, где живёт литва), и по форме.ЛитваПоверхность — равнина, со следами древнего оледенения. Поля и леса занимают 57 % территории, леса и кустарники — 30 %, болота — 6 %, внутренние воды — 1 %.ЛитваВысшая точка — 293,84 м над уровнем моря — холм Аукштайс (или Аукштасис калнас) в юго-восточной части страны. 23,5 км от Вильнюса.ЛитваКрупнейшие реки — Неман и Вилия.ЛитваБолее 3 тыс. озёр (1,5 % территории); крупнейшее из них — Друкия на границе Латвии, Литвы и Белоруссии (площадь 44,8 км²), самое глубокое — Таурагас, 61 м), самое длинное — Асвея длиной в 30 км у местечка Дубингяй.ЛитваКлимат переходный от морского к континентальному. Средняя температура зимой −5 °C, летом +17 °C. В среднем выпадает 748 мм осадков в год.ЛитваПолезные ископаемые: торф, минеральные материалы, строительные материалы.ЛитваТерритория современной Литвы была заселена людьми с конца X—IX тысячелетия до н. э. Жители занимались охотой и рыболовством, использовали лук и стрелы с кремёвыми наконечниками, скребки для обработки кожи, удилки и сети. В конце неолита (III—II тысячелетия до н. э.) на территории современной Литвы проникли индоевропейские племена. Они занимались земледелием и скотоводством, при этом охота и рыболовство оставались основными занятиями местных жителей вплоть до широкого распространения железных орудий труда. Индоевропейцы, заселившие земли между устьями Вислы и Западной Двины, выделились в отдельную группу, названную учёными балтами.ЛитваТрадиционно считается, что этническая основа Литвы сформирована носителями археологической культуры восточнолитовских курганов, сложившейся в V веке н. э. на территории современных Восточной Литвы и Северо-Западной Белоруссии. Около VII века литовский язык отделился от латышского.ЛитваСтановление государственности относится к XIII веку, при этом само название «Литва» впервые упомянуто в Кведлинбургских анналах под 1089 годом в сообщении об убийстве язычниками миссионера Бруно на границе Руси и Литвы (— косов.п.). По наиболее распространённой, но аргументировано опровергнутой, версии, топоним возник от названия небольшой реки Летаука, притока Нариса. Согласно более современной гипотезе, название страны могло произойти от этнонима «летяи» или «лейти», которым жители окрестных земель называли друминок литовских князей.ЛитваНачало XIII века в землях балтов-язычников с запада началось вторжение немецких рыцарей крестового похода. ЛитваЛитва, долгое время служившая основным предметом споров с крестоносцами.ЛитваВеликий князь Казимир, одновременно бывший и королём польским, расширил влияние династии Ягеллонов — подчинил Пруссию, послал своего сына на чешский и венгерский троны. В 1492—1526 годах существовала политическая система государства Ягеллонов, охватывавшая Польшу (с вассалами Пруссией и Молдавским княжеством), Великое княжество Литовское, Чехию и Венгрию.ЛитваПравовой основой государства являлся статут, изданный в трёх редакциях (1529, 1566, 1588), отражающих социально-экономические и политические изменения. Статут регламентировал вопросы гражданского, уголовного и процессуального права. На территории Великого княжества третья редакция статута действовала до 1840 года.Литва 1569 году в Люблине была заключена новая уния с Польшей, в результате которой образована Речь Посполитая. Согласно акту Люблинской унии литовый и польский правитель совместно избираемый королём, а государственные дела решались в общем сейме. Однако правовые системы, арми и чиновники оставались раздельными.Литва XVI—XVIII веках в Литве по польскому

Результат

 {{id: "7", "url": "https://ru.wikipedia.org/wiki?curid=7", "title": "Литва", "text": "ЛитваЛитва́ официальное название литовская республика государство расположенное северной ча сти европы площадь км² протяжённость севера юг 280 км запада восток 370 км население составляет человек январь 2021 занимает 137-е место мире численности населения 121-е территор ии имеет выход балтийскому морю расположена восточном побережье береговая линия составляет 99 км наименьший показатель среди государств балтии севере граничит латвией юго востоке белоруссией юго западе польшей калининградской областью россия площадь населения является самым крупным прибалтийским государством сто-ли-ца вильнюс официальный язык литовский язык де-юре не признан прибалтийским государством. Независимости страны провозглашено 11 марта 1990 года. 6 сентября 1991 года Государственный совет СССР признал независимость Литвы.ЛитваЛитва — член ООН (1991), ОБСЕ (1991), Совета Европы (1993), ВТО (2001), Европейского союза (2004), НАТО (2004) и ОЭСР (2018). Входит в Шенгенскую зону и Еврозону.ЛитваИтология слова «Литва» точно известна существует множество версий, ни одна из которых получила всеобщего признания. Корень «лит» варианты «лет/лёт» допускают различные толкования балтских славянских других индоевропейских языках например существуют созвучные топонимы территории словакии Литва румынии Литва известные XI-XII веков мнение е поспелова топоним образован древнего названия реки Летава (Lietavā) лить русское Летаука феодалное княжество землям которого протекала эта река временем заняло ведущее положение название распространено всё государство повести временных лет XII век упоминается этноним литва полностью совпадающий названием местности Литва в смысле (территория где живёт литва) форме поверхность равнина следы древнего оледенения поля и леса занимают 57 % территории, леса и кустарники — 30 %, болота — 6 %, внутренние воды — 1 %.ЛитваВысшая точка — 293,84 м над уровнем моря — холм Аукштайс (или Аукштасис калнас) в юго-восточной части страны. 23,5 км от Вильнюса.ЛитваКрупнейшие реки — Неман и Вилия.ЛитваБолее 3 тыс. озёр. 1,5 % территории. крупнейшее Друкия на границе Латвии Литвы и Белоруссии площадь 44,8 км². самое глубокое Таурагас, 61 м. самое длинное Асвея длиной 30 км. местечко Дубингяй. климат переходный морского континентальному средняя температура зимой −5 °C, летом +17 °C. выпадает 748 мм осадков. год. полезные ископаемые: торф минеральные материалы строительные материалы. ЛитваТерритория современной Литвы заселена людьми конца X-XI тысячелетия н.э. жители занимались охотой рыболовством использовали лук стрелы кремёвыми наконечниками скребки обработки кожи удилки сети конце неолита III-II тысячелетия н.э. жители занимались охотой рыболовством оставались основными занятиями местных жителей вплоть широкого распространения железных орудий труда индоевропейцы заселившие земли устьями вислы западной двины выделились отдельную группу названную учёными балтами. Традиционно считается этническая основа литвы сформирована носителями археологической культуры восточнолитовских курганов сложившейся в V веке н.э. на территории современных восточной литвы северо западной белоруссии около VII века литовский язык отделился латышского становление государственности относится XIII веку само название литва впервые упомянуто в кведлинбургских анналах 1089 годом в сообщении об убийстве язычниками миссионера Бруно на границе руси литвы — косов.п. наиболее распространённой, аргументировано опровергнутой версии топоним возник названия небольшой реки Летаука притока Нариса согласно современной гипотезе название страны могло произойти этнонима «летяи» или «лейти» которым жители окрестных земель называли друминок литовских князей. ЛитваНачало XIII века земли балтов язычников запада началось вторжение немецких рыцарей крестоносцев покорили Пруссию Ливонию Это время когда началась экспансия галицко волынского княжества середине XIII века многие литовские земли объединены князем миндовга принявшего 1251 году католичество. ЛитваНачало XIV века в землях балтов-язычников с запада началось вторжение немецких рыцарей крестоносцев покорили Литву. Литва, долгое время служившая основным предметом споров с крестоносцами. ЛитваВеликий князь Казимир, одновременно бывший королём польским расширил влияние династии Ягеллонов подчинил Пруссию послал своего сына чешский венгерский троны. 1492-1526 годах существовала политическая система государства Ягеллонов охватывавшая Польшу вассалами пруссией молдавским княжеством Великое княжество Литовское Чехию Венгрию. ЛитваПравовой основой государства являлся статут изданный трёх редакциях 1529-1566-1588 отражающих социально экономические политические изменения. Статут регламентировал вопросы гражданского уголовного процессуального права территории Великого княжества третья редакция статута действовала 1840 года. 1569 году в Люблине заключена новая уния польшей в результате которой образована Речь Посполитая. Согласно акту Люблинской унии литовый и польский правитель совместно избираемый королём, а государственные дела решались в общем сейме. Однако правовые системы, армия чиновники оставались раздельными. XVI-XVIII веках в Литве по польскому образцу сложилась политическая система известная шляхетская демократия характеризовалась наличием широких прав шляхты дворянства управлением государством одновременно этим происходила колонизация шляхты выраженная перенятием правшим образом великого княжества литовского польского языка культуры идентичности неприевилегированные сословия колонизация столь значительного влияния оказала XVIII веке в результате опустошительных войн всеобщего государственного кризиса Речь Посполитая пришла упадку подала влияние российской империи 1772-1793-1795 годах состоялись разделы Речи Посполитой Россией Пруссией Австрией вся территория бывшего великого княжества литовского присоединена российской империи попытках восстановить государственность польское литовское дворянство приняло сторону Наполеона 1812 году также неоднократно поднимало восстания 1830-1831-1863-1864 годах однако окончилась поражением стремлении ликвидировать польское влияние в Литве российский власти предприняли широкую кампанию деполонизации русификации 1864 году частично запрещена литовская печать латиницей литовское население особенно католическое духовенство с сопротивлялись русификации кириллические издания игнорировали книги напечатанные латиницей книгошопы нелегально ввозили соседней пруссии 1904 году запрет литовскую латиницу отменён и начавшаяся первая мировая война быстро распространилась территории литвы концу 1915 года этнические литовские земли контролировались германией литовцы потеряли политические права в начале запрещены литовские периодические издания однако литовская интеллигенция попыталась воспользоваться геополитической ситуацией начала искать возможности восстановления независимости литвы 19-22 сентября 1917 года вильнюсе проведена литовская конференция время которой избрана литовская триба совет литвы ходе конференции принято решение необходимости создания независимого литовского государства этнографических границах столицей вильнюсе председателем совета избран сиемона 11 декабря 1917 года ещё 16 февраля 1918 года провозглашено восстановление литовского государства отличие принятой 11 декабря литовку германских властей декларация документ 16 февраля говорит полной независимости литвы россия германия однако документ 16 февраля даёт самостоятельность бумаге заключения брест литовского мирного договора германия игнорирует декларацию 16 февраля ссылался резолюцию 11 декабря взвешивает возможность создать литовское корольство германским монархом 23 марта 1918 года император вильгельм II признал независимость литвы основании акта признании литов

Как видно из результата токенизации, не все токены были выделены удачно, к примеру ‘ос’, ‘iii’, ‘→’, знаки ударения все еще присутствуют. Добиться лучших результатов можно более тщательной обработкой текста.

Из достоинств выбранного метода токенизации можно выделить:

- использование готового словаря стоп-слов русского языка

Из недостатков:

- невысокая скорость обработки

Исходный код

```
1  import os
2  import json
3  from pathlib import Path
4  from trash import trash
5  import nltk
6  from nltk.corpus import stopwords
7  from nltk.tokenize import word_tokenize
8  nltk.download("stopwords")
9
10 russian_stopwords = stopwords.words("russian")
11
12 def replacer(text, dic):
13     for i, j in dic.items():
14         text = text.replace(i, j)
15     return text
16
17
18 def preprocess_text(text):
19     text = text.lower()
20     text = replacer(text, trash)
21
22     tokens = word_tokenize(text, language="russian")
23     tokens = [token for token in tokens if token not in russian_stopwords]
24
25     text = " ".join(tokens)
26     text = " ".join(text.split())
27     return text
28
29
30 if __name__ == '__main__':
31
32     raw_dir_path = "/Users/denis/MAI/IR/lab3/wikipedia_articles_one_by_one"
33     raw_dir_list = [(raw_dir_path + "/" + i) for i in sorted(
34         os.listdir(raw_dir_path),
35         key=lambda x: int(x.replace(".json", "").replace("wiki_article_", ""))
36     )]
37     new_path = "/Users/denis/MAI/IR/lab3/wikipedia_articles_tokenized"
38
39     for filepath in raw_dir_list:
40         with open(filepath) as file:
41             data = json.load(file)
42             data["text"] = preprocess_text(data["text"])
43             new_name = f"wiki_article_{data['id']}.json"
44             with open(new_path + "/" + new_name, "w", encoding='utf8') as new_file:
45                 json.dump(data, new_file, ensure_ascii=False)
46
```

Выводы

В процессе выполнения данной лабораторной работы была выполнена токенизация текстов статей википедии с использованием библиотеки nltk на Python. Было интересно поискать готовые средства для предобработки текста и создания токенов. Наличие встроенного и постоянно обновляемого словаря стоп-слов в библиотеке nltk было приятным открытием.