

به نام خدا

امیرحسن سعادت‌مند

گزارش فاز یک پروژه بازیابی

۹۷۳۱۰۲۹

(۱)

طبق مراحل گفته شده در دستور کار و مطالب تدریس شده در کلاس چهار گام اصلی انجام میشود

(۱) نرمالایز کردن هر داکيومنت

(۲) توکنایز کردن (استخراج توکن از هر سند)

(۳) ریشه یابی (steming)

(۴) حذف stopwords ها (کلمات پرتکرار)

ابتدا data structure ها و کلاس های ساخته شده را بررسی میکنیم

```
import os

import pandas as pd

class Adapter:
    base_dir = os.path.dirname(__file__)
    file_name = "IR1_7k_news.xlsx"
    data_path = os.path.join(base_dir, file_name)

    def read_excel_file(self):
        df = pd.read_excel(self.data_path)
        return df

    def get_doc_id_and_documents(self):
        df = self.read_excel_file()
        for doc_id, document in enumerate(df.values):
            yield doc_id, document[0], document[2]

get_doc_id_and_document = Adapter().get_doc_id_and_documents()
```

کلاس Adapter وظیفه خواندن از فایل اکسل به صورت داک به داک (تکی تکی) و به صورت جنریک نوشته شده تا قابل iterate داشته باشد و نیاز نباشد در وهله اول همه چیز خوانده شود

کلاس MyTokenizer است که در آن Stopword ها را از فایلی به نام stopwords.bat در پروژه هضم است دربردارد به علاوه آبجکت iterable از داکيومنت هایی که از adapter فرستاده میشوند

متد اصلی این کلاس get_tokenized_stream است .. به طوری نوشته شده که هرداک یکی یکی از Adapter گرفته شده و بعد از نرمالسازی و حذف stopword ها و ریشه یابی استریمی از توکن برمیگرداند
return Tupe(tokens , doc_id)

```

from parsivar import Normalizer, Tokenizer, FindStems

class MyTokenizer:
    def __init__(self):
        self.stop_words = self.get_stop_words()
        self.docs = get_doc_id_and_document
        self.stemer = FindStems().convert_to_stem

    def get_tokenized_stream(self) -> Tuple:
        for doc_id, document, doc_title in self.docs:
            normalized_document: str = Normalizer().normalize(document)
            token_words: List = Tokenizer().tokenize_words(normalized_document)
            tokens = list()
            for word in token_words:
                stem = self.stemer(word)
                # if stem not in self.stop_words:
                tokens.append(stem)
            yield doc_id, tokens, doc_title

    @staticmethod
    def get_stop_words() -> List:
        base_dir = os.path.dirname(__file__)
        file_name = "stopwords.dat"
        with open(os.path.join(base_dir, file_name), 'r') as file_p:
            stop_word_list = file_p.readlines()
            stop_word_list = [stopword.replace("\n", " ") for stopword in stop_word_list]
        return stop_word_list

get_token_stream = MyTokenizer().get_tokenized_stream()

```

دیتا اسراکچر اصلی برای شاخص از نمونه تمرین hw1 موجود در صفحه درس نمونه برداری شده

```

from typing import Dict

class DocTerm:
    def __init__(self, doc_id, term_id):
        self.count = 0
        self.docId = doc_id
        self.termId = term_id
        self.positions = list()

    def insert(self, position):
        self.count += 1
        self.positions.append(position)
        self.positions = sorted(self.positions)

class Term:
    def __init__(self, id):
        self.count = 0
        self.id = id
        self.docTerms: Dict[DocTerm] = dict()

    def insert(self, doc_id, position):
        self.count += 1
        if doc_id not in self.docTerms.keys():
            self.docTerms[doc_id] = DocTerm(doc_id, self.id)
            sorted(self.docTerms.items())
        self.docTerms[doc_id].insert(position=position)

```

علت اصلی نرمالایز کردن یکسان سازی نوشته هاست به عنوان مثال
U_S_A , U.S.A به انگلیسی باشند یا چندین فرم از یک کلمه به یک صورت نوشته شوند مثل

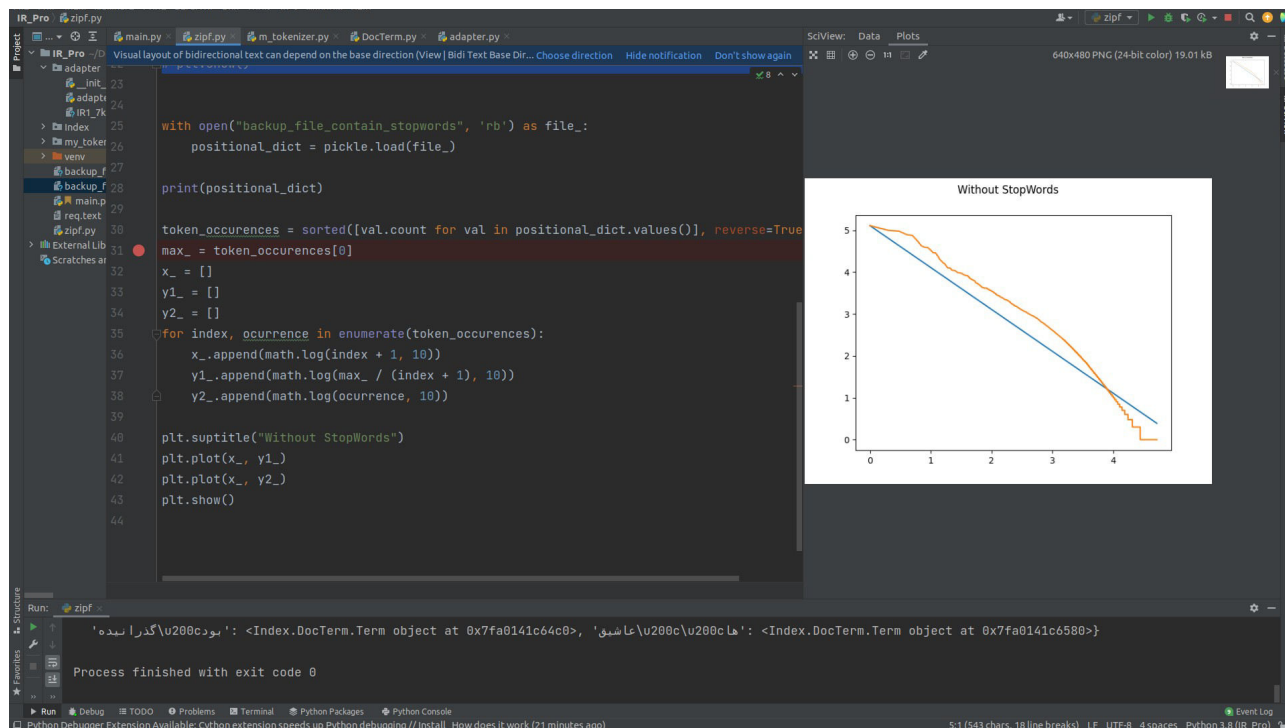
با tokenize کردن تک تک کلمات استخراج میشوند و آماده ساخت دیکشنری کلمات میشوند

و با ریشه یابی حالت های مختلف یک مفهوم و یا یک فعل را به یک صورت مینویسیم
مثلا همه افعال میگویم بگویم و ... از گفتن هستند

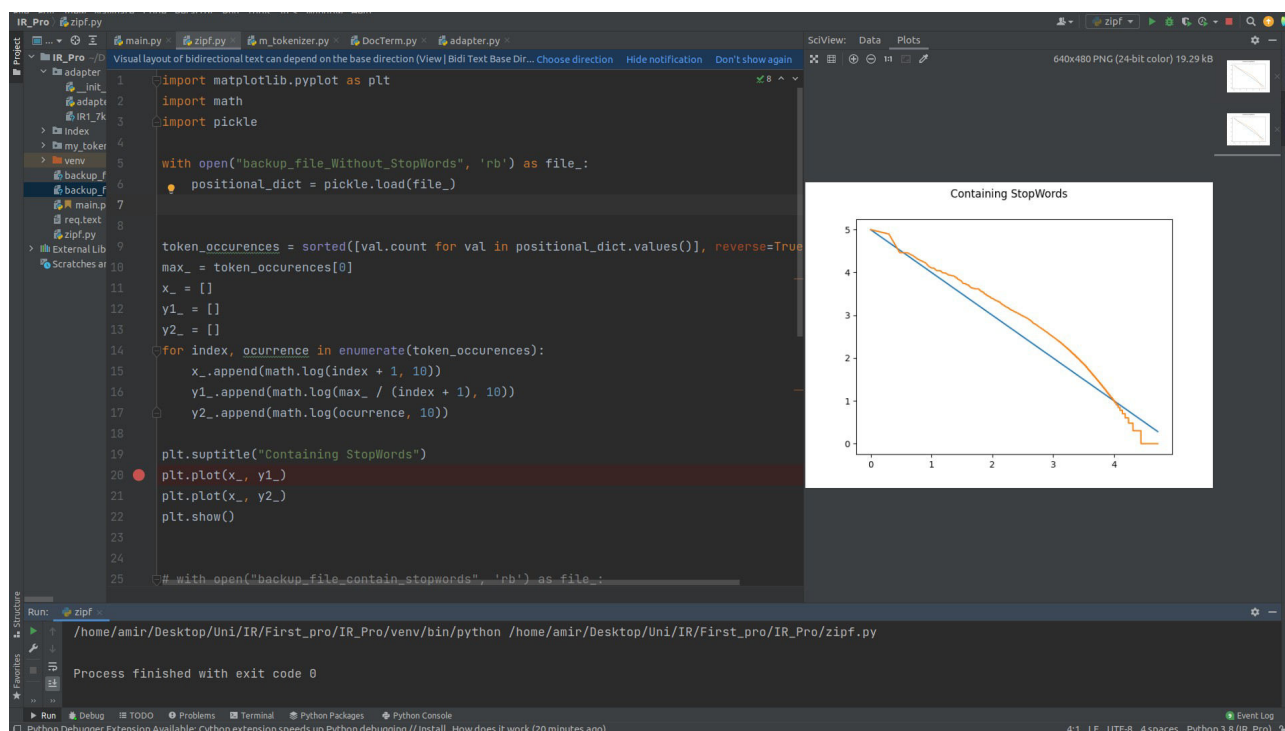
و در نهایت عمل حذف stopwords هاست که کلماتی همچون حروف اضافه و ... که خاصیت مفهومی برای بازایی ندارند را شاخص نمیکیم

(۲

نمودار پس از حذف stopwords ها :



و با در نظر گرفتن آنها (شامل stopwords ها) :



(3)

قبل ریشه یابی :

تعداد توکن ها ۳۲۶۳۹۲

تعداد کلمه ها ۱۷۹۰۴

تعداد همه توکن ها ۱۹۵۹۵۵۲

تعداد کل کلکلمات ۷۲۸۹۰

این در ازای $b = 0.67$ به دست میاد

بعد ریشه یابی :

تعداد توکن ها ۳۳۹۸۸۳

تعداد کلمه ها ۱۳۸۱۲

تعداد کل توکن ۲۰۷۸۶۸۴

تعداد کل کلم ها ۵۵۸۹۰

و در این حالت هم $b = 0.66$ به دست میاد و رابطه برقرار هست

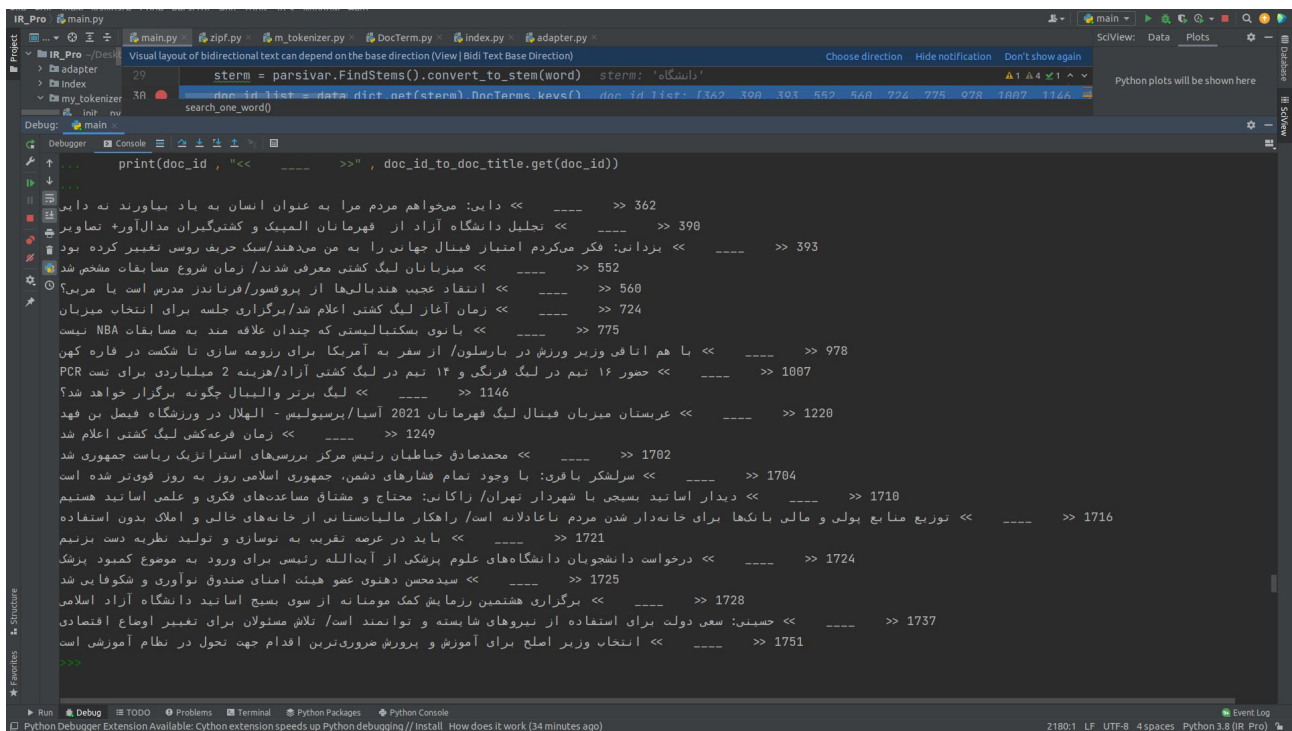
(4)

کلماتی که در ریشه خودشان پسوند هایی دارند که بعد از ریشه یابی ب مشکل میخورند و شکل درست آنها ذخیره نمیشود

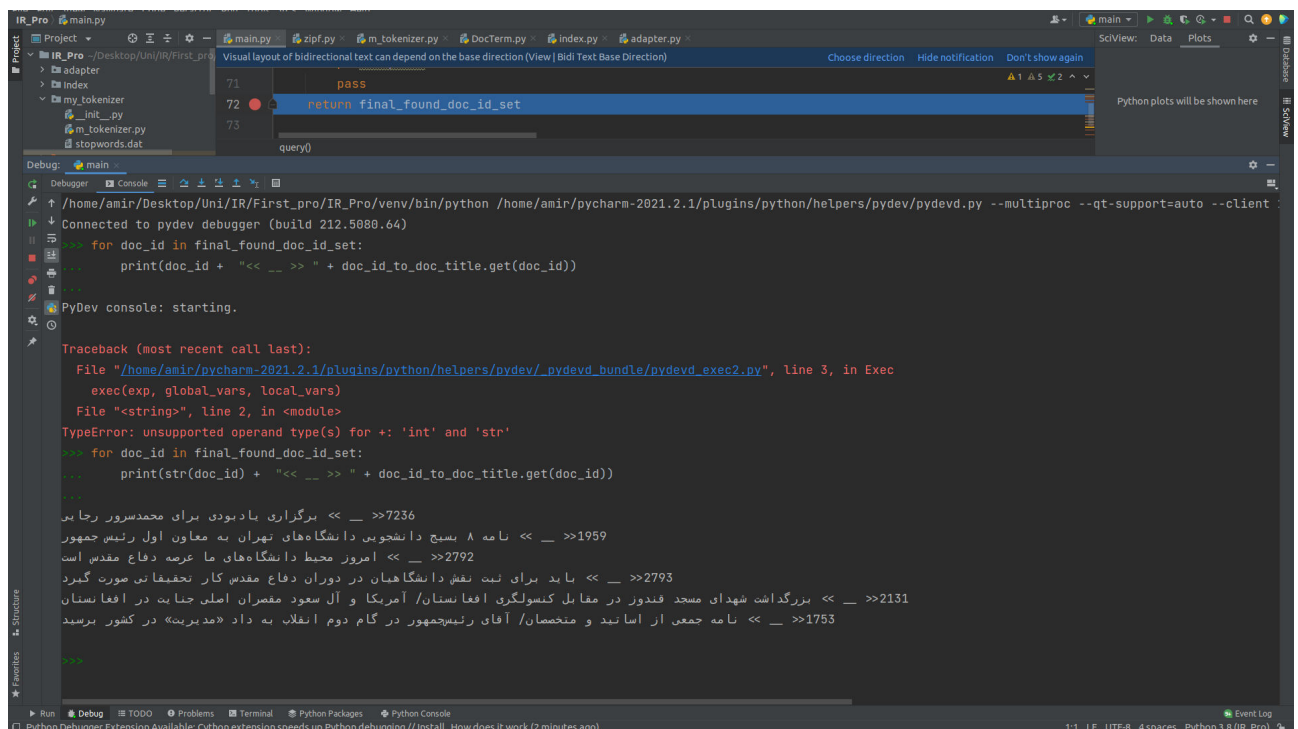
کلمه ای مانند درخت که به ت ختم شده اما بعد ریشه یابی به صورت درخ ذخیره شده
یا فعلی مثل نهادن یه هر وجه و زمانی از این مصدر که به نه ختم میشود و با ۹ اشتباه گرفته میشود

(5)

الف (پرسمان = دانشگاه

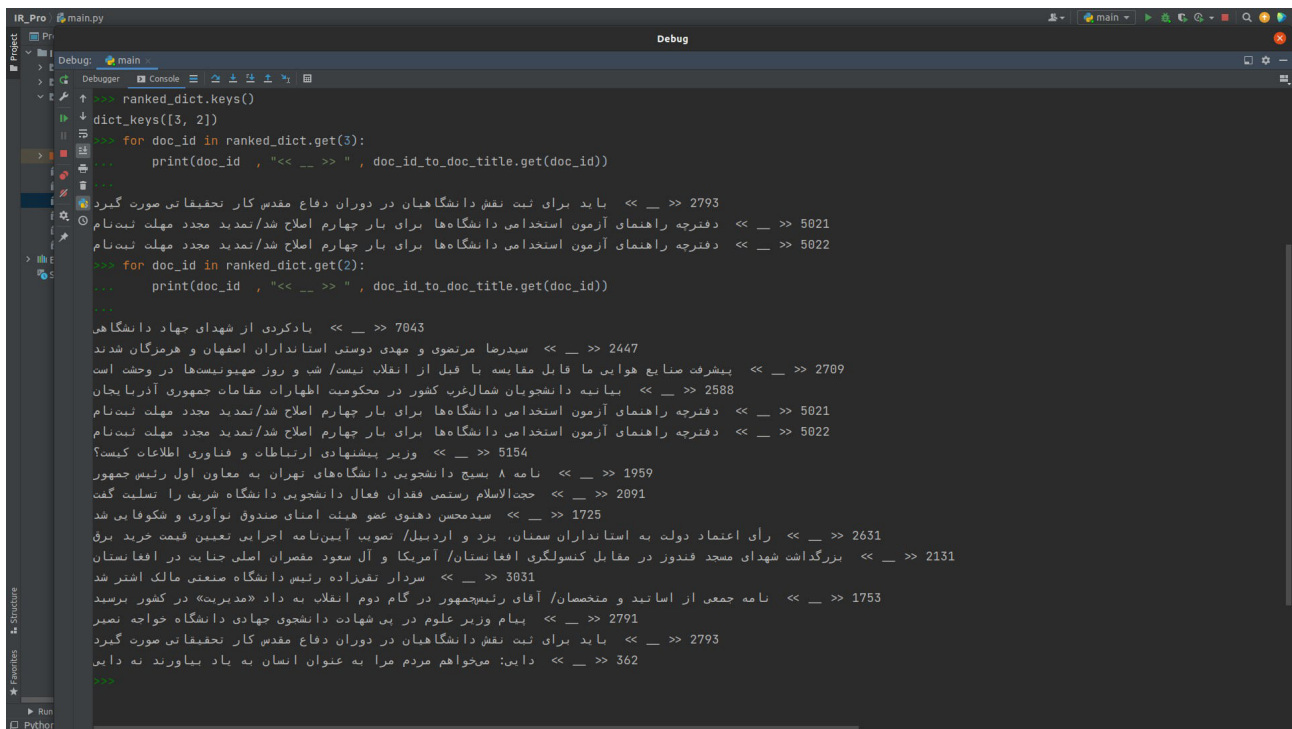


ب) دانشگاه امیرکبیر

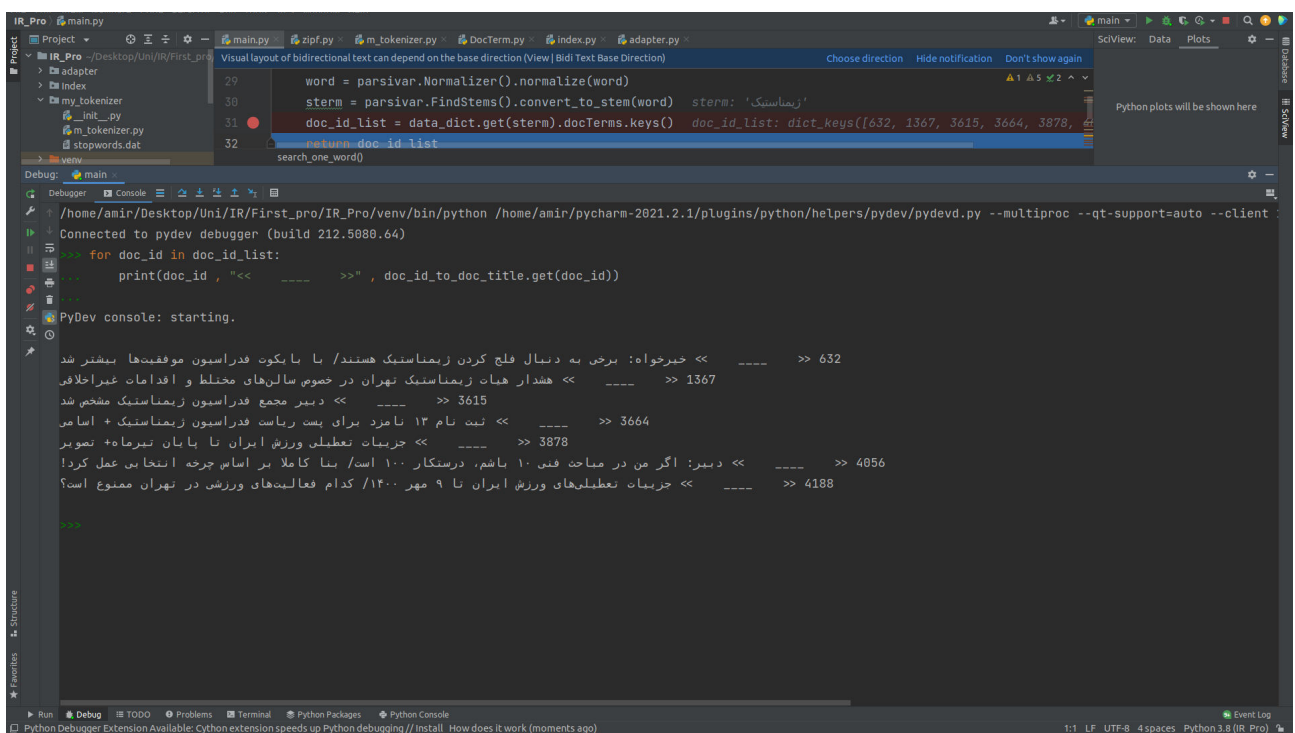


پ) دانشگاه صنعتی امیرکبیر

در اینجا داده ها رنگینک شدند و به ترتیب بالاترین رنک را دارند چاپ شده اند



ت (پرسمان زمیناستیک



ث) واکسن استرازنکا

