

((گزارش فاز اول پروژه بازیابی اطلاعات))

۹۷۳۱۱۲۲

پویان حسابی

(۱) با ذکر مثال شرح دهید که در گام پیش پردازش چه عملیاتی انجام داده‌اید. همچنین دلیل انجام هر پردازش را ذکر کنید.

در گام پیش پردازش مطابق با دستور پروژه و مباحث درس ابتدا رشته مربوطه را نرمال می‌کنیم برای اینکه یکسری از کاراکترها و به طور مثال اعداد فرم یکسانی بگیرند، سپس توکن‌های آن را استخراج می‌کنیم (هر کلمه به صورت جداگانه) و کلماتی از هر content (محتوای سند) را داریم. از آنجایی که نحوه ریشه‌یابی (stemming) بر اساس لغات در کتابخانه مربوطه انجام می‌شود سپس این تابع را فراخوانی می‌کنیم، دلیل آن این است که کلماتی مثل رفتیم، رفتند و ... یک کلمه حساب شوند و به نوعی به یک فرم باشند. سپس چک می‌کنیم که کلمه مربوطه stop word نباشد چرا که تاثیری در نتایج ما ندارد و صرفاً عملیات را زیاد می‌کند، نمونه‌ای از استاپ ورد: از، به، تا، و ...

از آنجایی که در پایتون دسترسی به مقادیر دیکشنری (ساختمان داده dict) در $O(1)$ انجام می‌شود دسترسی به محتوای هر سند را بر اساس $\{url(str) : content(list\ of\ tokens)\}$ در نظر گرفته شده است و از طرفی محتوای ما که شامل توکن‌هایی از کلمات است باید در لیست قرار گیرد تا مکان آنها برای positional index استفاده شود.

```
def pre_processing():
    start_time = time.time()
    tokenized_dict = {}
    stopwords = stopwords_list()
    stemmed_text = []
    for url in url_content:
        tokenized_dict[url] = word_tokenize((Normalizer().normalize(url_content[url])))
        for word in tokenized_dict[url]:
            if word not in stopwords:
                stemmed_text.append(Stemmer().stem(word))
        tokenized_dict[url] = stemmed_text.copy()
        stemmed_text.clear()
    print("Time spent for pre processing: " + str(time.time() - start_time))
    return tokenized_dict
```

از کتابخانه hazm استفاده شده است که نمونه هایی از هر کدام از مراحل در پایین آورده شده:

```
>>> normalizer = Normalizer()
>>> normalizer.normalize('اصلاح نویسه ها و استفاده از نیمفاصله یردازش را آسان می کند')
'اصلاح نویسه ها و استفاده از نیمفاصله یردازش را آسان می کند'
```

```
>>> word_tokenize('ولی برای یردازش، جدا بهتر نیست؟')
['ولی', 'و', 'برای', 'یردازش', 'و', 'جدا', 'و', 'بهتر', 'و', 'نیست', 'و', '؟']
```

```
>>> stemmer = Stemmer()
>>> stemmer.stem('کتابها')
'کتاب'
```

```
print(Stemmer().stem("کتابها"))
print(Stemmer().stem("کتابها"))
```

```
Test
"C:\Program Files\Python\python.exe"
```

```
Test
"C:\Program Files\Python\python.exe"
```

```
print(stopwords_list())
```

```
Test
"C:\Program Files\Python\python.exe" "D:/Edu/Term 8/Information Retrieval/Project/Final IR Project/Test.py"
```

```
Test
"بار" و "ابتدا" و "حق" و "الفاظ" و "می" و "200C" و "طایقی" و "هون" و "نشد" و "یافت" و "پر" و "اطرفشان" و "کد" و "جسمی" و "ادب" و "نوازه" و "می" و "200C" و "اذا" و "رله" و "جهد" و "انجا" و "ار"
```

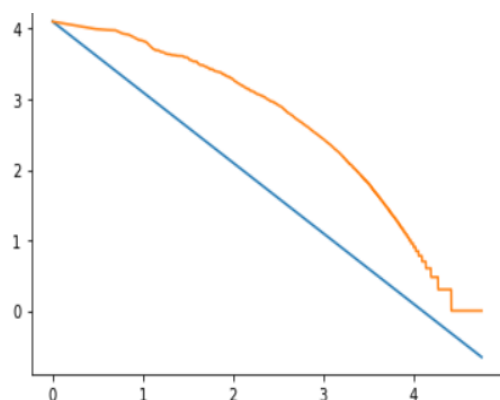
۲) صحت قانون Zipf را در دو حالت قبل و بعد از حذف کلمات پر تکرار از واژه نامه بررسی کنید(رسم نمودار برای هر حالت الزامی است). در صورت برقراری / عدم برقراری این قانون در هر حالت، علت را شرح دهید.

Zipf's law: The i^{th} most frequent term has frequency proportional to $1/i$.

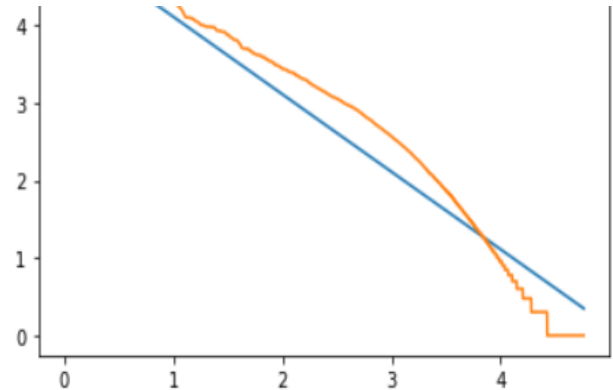
$cf_i \propto 1/i = K/i$ where K is a normalizing constant

cf_i is collection frequency: the number of occurrences of the term t_i in the collection.

حذف stop words:



حذف نکردن stop words:



همانطور که می‌دانیم و گفته شد قانون zipf مبتنی بر تعداد تکرار کلمات است. در واقع رابطه عکس با کلمات پر تکرار دارد. برای stop words از کتابخانه hazm استفاده شده که دقت کمتری نسبت به کتابخانه دیگر دارد ولی سرعت بسیار بیشتری دارد. همانطور که مشاهده می‌کنید با حذف کلمات پرتکرار فاصله زیادی با تخمین ایده آل (خط آبی) دارد ولی با حذف نکردن کلمات پر تکرار نمودار نارنجی خیلی نزدیکتر به نمودار آبی می‌باشد پس تخمین ما در حالت با بدون حذف stop words براساس قانون zipf بوده ولی وقتی آنها را حذف کنیم فاصله آن زیاد می‌شود. و در تئوری هم همچنین نتیجه ای انتظار می‌رفت چون با حذف کلمات پرتکرار در واقع آنها را از نمودار حذف می‌کنیم و نمودار نزدیک به داده های دیگر می‌شود و پراکندگی زیاد می‌شود.

۳) صحت قانون heaps:

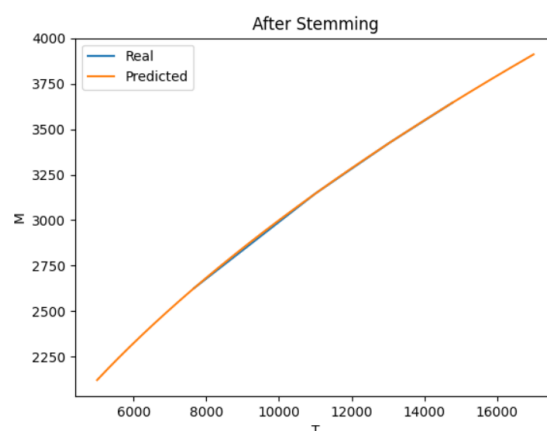
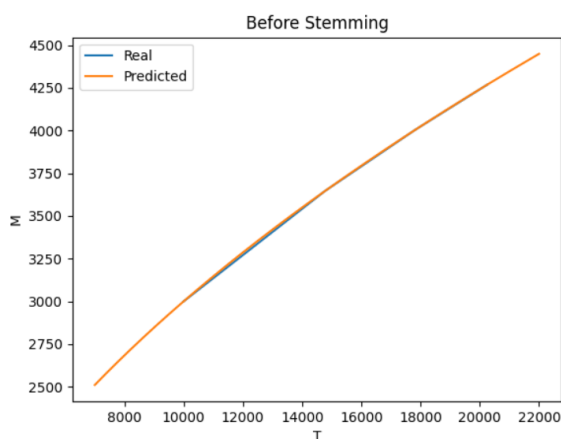
قبل از ریشه یابی:

بعد از ریشه یابی:

تعداد اسناد	تعداد کلمات مجزا	تعداد کل توکن ها	تعداد اسناد	تعداد کلمات مجزا	تعداد کل توکن ها
۵۰۰	۹۹۸۸	۷۵۵۷۸	۵۰۰	۷۶۶۵	۴۷۶۲۹
۱۰۰۰	۱۴۷۶۸	۱۵۴۳۵۵	۱۰۰۰	۱۱۰۰۶	۹۷۰۴۵
۱۵۰۰	۱۷۸۲۱	۲۳۲۹۷۹	۱۵۰۰	۱۳۰۴۷	۱۴۵۶۱۳
۲۰۰۰	۲۰۲۵۰	۳۰۶۰۳۲	۲۰۰۰	۱۴۷۶۶	۱۹۱۴۷۲

$$M = k * \text{power}(T, b)$$

. b = 0.5 و k = 30 و T: number of tokens و M: size of vocabs



بر اساس نمودار ها مشاهده می شود تقریبا هر دو از قانون heaps تبعیت می کنند. در حالتی که بعد از ریشه یابی است مقدار نزدیک تری به مقداری که پیشبینی شده است دارد و به این معناست که ریشه یابی باعث می شود که قانون heaps پر رنگ تر شود و عملا بعد از آن شبیه تر است، در نتیجه می توان مشاهده کرد که با حذف کلمات هم ریشه علاوه بر اینکه حجم لغات کاهش می یابد دقت مدل ما نیز افزایش می یابد و مدل به حالت واقعی نزدیک می شود. و احتمالا هر چه تعداد اسناد بیشتر شود، بیشتر به این مقدار نزدیک می شود که به این معناست ریشه یابی لغات بسیار سودمند می باشد.

۴) حداقل سه مورد از مواردی که در ریشه یابی با چالش روبرو بودید را ذکر کنید

فوتبال امی رفتند رفت رفی انسی گلاب

```
array = ["فوتبال", "است", "روشن", "رحیم", "رحم", "انسان", "کباب"]
for word in array:
    print(Stemmer().stem(word), end=" ")
```

مشاهده می‌شود کلمه (است) به (اس) تبدیل شده و به طور کلی حروف بعد از آن را حذف کرده است. کلمه (رفتند) همان (رفتند) باقی مانده و بدون تغییر است، کلمه (رفتیم) به درستی بوده ولی کلمه (رفتم) تبدیل به (رف) شده است و کلمه (انسان) به اشتباه ریشه یابی شده است.

۵) پاسخ به پرسمان در حالت‌های زیر:

فرم کلی پاسخ به کوئری به این شکل می باشد که ابتدا شماره سند، url، عنوان چاپ می شود.

```
def show_ranked_documents(ranked_docs: list):
    for counter, tuple_ in enumerate(ranked_docs):
        if counter > 4:
            break
        doc_id = tuple_[0]
        print("***\nResult " + str(counter + 1))
        print("Document ID: " + str(doc_id))
        print("URL: " + str(urls[doc_id]))
        print("Title: " + str(url_title[urls[doc_id]]))
```

الف) تحریم های آمریکا علیه ایران

Result 1

Document ID: 9418

URL: <https://www.farsnews.ir/news/14001005000843/> بازخوانی - بیانات رهبر انقلاب - در دیدار دست‌اندرکاران مراسم سالگرد

Title: بازخوانی | بیانات رهبر انقلاب در دیدار دست‌اندرکاران مراسم سالگرد شهادت سردار سلیمانی

رفع تحریم دست دشمن است. ژنرال‌های جنگ نرم آمریکا و استکبار هستند. با ملت ایران بدی کرد. بظاهر مذهبی متمایل به یکی از فرقه‌ها و علیه مقاومت را ایشان مدتها قبل پیش‌بینی کرد

تحلیل: مشاهده می‌شود تنها یک سند بازگردانی شد چرا که اسنادی که هم آمریکا و هم ایران و هم علیه داشته باشند بسیار کم می‌باشند. این سند تا حدی مرتبط به کوئری کاربر است از این جهت که علیه را به عنوان یک استاپ ورد در نظر نگرفته و تاثیر بسزایی در کوئری گذاشته است. از طرفی لغات تحریم و آمریکا و ایران همگی در سند مربوطه موجود است.

ب) تحریم های آمریکا ! ایران

Result 1

Document ID: 6929

URL: <https://www.farsnews.ir/news/14001222000450/> توضیحات - یک منبع - آگاه - درباره - وقفه - مذاکرات - وین

Title: توضیحات یک منبع آگاه درباره وقفه مذاکرات وین

برداشتن تحریم‌ها وعده داده عمل کند و تضمین معتبری هم برای انجام این تعهدات داده شود. آمریکا در مقابل عمل به این تعهد برجای مقاومت می‌کند

Result 2

Document ID: 7209

URL: <https://www.farsnews.ir/news/14001214001141/> بحران - اوکراین - مصداق - ماهیت - مافیایی - آمریکاست

Title: بحران اوکراین، مصداق ماهیت مافیایی آمریکاست

Result 3

Document ID: 7183

URL: <https://www.farsnews.ir/news/14001215000516/> میزگرد - تلاطم - در - اوکراین | ائتلاف - اقتصادی - حول - محور - دلار - فرو - می‌باشد

Title: میزگرد «تلاطم در اوکراین» | ائتلاف اقتصادی حول محور دلار فرو می‌باشد

Result 4

Document ID: 9742

URL: <https://www.farsnews.ir/news/14000924000773/> مرکز پژوهش‌های مجلس-مذاکرات-وین-به-توافقی-زودهنگام-منجر-نمی‌شود

Title: مرکز پژوهش‌های مجلس: مذاکرات وین به توافقی زودهنگام منجر نمی‌شود

Result 5

Document ID: 11864

URL: <https://www.farsnews.ir/news/14000803000676/> اهرم‌سازی-از-افغانستان-در-برجام-نقطه-عزیمت-آمریکا-در-مذاکرات-جامع-با

Title: اهرم‌سازی از افغانستان در برجام/ نقطه عزیمت آمریکا در مذاکرات جامع با ایران چیست؟

تحلیل: نتایج تعداد زیادی از اسناد می‌باشد که صرفاً ۵ تا از آنها نمایش داده شده است. همانطور که دیده می‌شود هر ۵ سند کاملاً مرتبط به کوئری می‌باشد و همگی مربوط به تحریم‌های آمریکا می‌باشد ولی یک باگ اساسی دارد این کوئری آن هم لغت ایران می‌باشد که وقتی ریشه‌یابی می‌شود کلمه ایرا می‌شود که به این شکل کوئری و نتایج را خراب می‌کند.

پ) کنگره ضد تروریست

Result 1

Document ID: 6929

URL: <https://www.farsnews.ir/news/14001222000450/> توضیحات-یک-منبع-آگاه-درباره-وقفه-مذاکرات-وین

Title: توضیحات یک منبع آگاه درباره وقفه مذاکرات وین

ریگان که نمی‌خواست کمتر از **کنگره ضد تروریست** جلوه کند

تحلیل: کلمه «ضد تروریست» تنها در یک سند آمده است. در اسناد دیگر کلمه ضد و تروریست جداگانه آمده که باعث می‌شود آنها در نتایج نیابند. اینجاست که مفهوم spell correction اهمیت پیدا می‌کند چرا که کاربر کلمه را به اشتباه وارد کرده است.

ت) "تحریم هسته‌ای" آمریکا!! ایران

Result 1

Document ID: 9496

URL: <https://www.farsnews.ir/news/14000926000385/> گفت‌وگوی-مشروح-| -ترقی-آمریکا-شروط-ایران-را-نپذیرد-پشت-در-مذاکرات

Title: گفت‌وگوی مشروح | ترقی: آمریکا شروط ایران را نپذیرد، پشت در مذاکرات می‌ماند/ روحانی کشور را به بن‌بست کشاند

محدود کردن مذاکره به موضوع رفع یکباره تحریم‌ها است. که موضوعی در مورد هسته‌ای باقی نمانده که بخواهد در مورد آن مذاکره کند. سه شرط را برای بازگشت آمریکا به برجام مطرح کرده است

Result 2

Document ID: 7230

URL: <https://www.farsnews.ir/news/14001214000432/> رئیس-دانشکده-هسته‌ای-دانشگاه-شهید-بهشتی-انرژی-هسته‌ای-نیاز-قطعی-امروز

Title: رئیس دانشکده هسته‌ای دانشگاه شهید بهشتی: انرژی هسته‌ای نیاز قطعی امروز و فردای ایران است

Result 3

Document ID: 11495

URL: <https://www.farsnews.ir/news/14000813000113/> گزارش-فارس-از-اجتماع-میدانی-یوم-الله-۳-آبان-در-تهران-تاکید-بر-ضرورت

Title: گزارش فارس از اجتماع میدانی یوم الله ۱۳ آبان در تهران/ تاکید بر ضرورت هوشیاری ملت‌های منطقه برابر توطئه‌های آمریکا

Result 4

Document ID: 11341

URL: <https://www.farsnews.ir/news/14000816000568/> ماجرای-دستگیری-دو-دانشمند-ایرانی-در-آمریکا-سلیمانی-مرا-در-بند

Title: ماجرای دستگیری دو دانشمند ایرانی در آمریکا/ سلیمانی: مرا در بند قاتلین و قاچاقچیان بین‌المللی زندانی کردند

Result 5

Document ID: 11489

URL: <https://www.farsnews.ir/news/14000813000157/> آمریکایی‌ها-عادت-کرده‌اند-مدام-از-ایران-شکست-بخورند-روایت-سرلشکر

Title: آمریکایی‌ها عادت کرده‌اند مدام از ایران شکست بخورند/ روایت سرلشکر سلامی از مواجهه سپاه با آمریکا در دریای عمان

تحلیل: این کوئری علی‌رغم آنچه در صورت پروژه گفته پیچیدگی زیادی ندارد چرا که Double quotation به راحتی از آن حذف می‌شود و تعداد کلمات بالا در کوئری سیستم ما را کند نمی‌کند زیرا فارغ از کوئری ماتریس‌ها ساخته می‌شوند. ۵ نتیجه‌ای که آمده همگی مرتبط هستند ولی همچنان مشکل کلمه ایران وجود دارد.

ث) اورشلیم ! صهیونیست

No result for given query, because of this word: اورشلیم

تحلیل: هیچ سندی کلمه اورشلیم را دارا نمی‌باشد برای همین هیچ نتیجه‌ای خروجی نداده است.