

وزارت علوم، تحقیقات و فناوری
دانشگاه تحصیلات تکمیلی علوم پایه
گاوزنگ، زنجان



پروژه‌ی درس پردازش زبان طبیعی

سارینا دانایی

۹۸۴۱۳۱

استاد: دکتر ابراهیم انصاری

پاییز ۹۸

تعریف مسئله:

در این پروژه قصد داریم که جست و جوی میان ۱۴۰۰ داکيومنت را انجام دهيم . به صورتی که برای هر query که در یک فایل جداگانه موجود میباشد . بهترین جوابها که در اینجا به معنای مرتبط ترین داکيومنت ها می باشد را برگردانيم.

که باید این کار را با دو روش که یک روش معمول و یک روش گفته شده در مقاله ی & salton backley میباشد انجام دهيم و خروجی را در نهایت مقایسه کنیم.

توضیح کد:

در بخش ۱- فایل cran که شامل ۱۴۰۰ داکيومنت می باشد را باز میکنيم و mode را برابر r قرار می دهيم که به معنای read یا خواندن می باشد و در این فایل نمیتوانيم عمل نوشتن انجام دهيم.

سپس فایل را در یک متغیر به نام token قرار می دهيم که یک لیست می باشد و هر خانه ی جدول شامل یک کلمه ی داکيومنت می باشد. (با دستور split کلمات داکيومنت را براساس space از هم جدا کردیم و در لیست قرار دادیم)

در بخش ۲- با استفاده از تابع FreqDist که یکی از توابع کتابخانه nltk می باشد تعداد هر کلمه در لیست token را چاپ می کنیم (با استفاده از دستور freq.most_common() به صورت نزولی چاپ می شوند)

در بخش ۳ با دستور `porterStemmer` تمامی کلمات `stem` می‌شوند به این معنا که تمامی پسوند های کلمات مثل `ing, ed, es, ...` از کلمات حذف میشوند.

نکته ۱: در این دستور از یک حلقه استفاده شده است که ۵ بار متوالی کلمات `stem` شوند. به این دلیل که در یک بار استفاده از این دستور یک سری از کلمات `stem` نشدند و ۵ بار دستور را اجرا کردم که جواب بهتری در نهایت به دست آمد اما زمان اجرا افزایش یافت.

در بخش ۴ لیست `stopword` ها چاپ شده که شامل ۱۷۹ کلمه می‌باشد.

سپس تمامی `stopword` ها با استفاده از دستور `remove` از داکيومنت ها حذف شدند. که در این دستور نیز به همان دلیل نکته ۱ پنج بار دستور `remove` انجام شده.

در متغیر `tag` یک لیست تعریف شده که تگ هایی که در داکيومنت هستند را در لیست وارد کردیم. این تگ ها نیز با استفاده از دستور `remove` حذف شدند.

در بخش ۵ می‌خواهیم که این ۱۴۰۰ داکيومنت را از هم جدا کنیم و دوباره در یک لیست می‌گذاریم به صورتی که بتوانیم با استفاده از ایندکس که همان شماره داکيومنت میباشد. داکيومنت ها را به صورت جداگانه بخوانیم. در این قسمت `split` را براساس `I`. انجام میدهیم چون داکيومنت ها براساس `I`. از هم جدا شده اند.

در بخش ۶ با استفاده از `pandas` یک `DataFrame` تعریف کردیم با نام `tf` که ستون های آن از ۱ تا ۱۴۰۰ (شماره ی داکيومنت ها) و سطرها ی آن شامل `type` ها (کل کلمات داخل داکيومنت میباشد)

سپس با استفاده از حلقه تعداد تکرار هر کلمه را در هر داکيومنت به دست آوردیم و برطول لیست شامل داکيومنت ها تقسیم کردیم تا مقدار `tf` را به دست آوریم.

سپس دوباره یک `DataFrame` دیگر با نام `idf` می‌سازیم که شامل ۲ ستون میباشد :

ستون اول باید مقدار `idf` در آن قرار گیرد با نام `idf`

و ستون دوم که با نام counter مشخص شده است تعداد تکرار هر کلمه که در سطر ها قرار دارد در کل ۱۴۰۰ داکيومنت نشان میدهد (یعنی هر کلمه در کل داکيومنت چند بار تکرار شده است) برای مقدار idf هر کلمه باید تعداد داکيومنت ها را بر تعداد تکرار آن کلمه در کل داکيومنت ها که همان ستون counter میباشد. به دست آوریم و سپس لگاریتم میگیریم.

در بخش ۷. دقیقاً مانند بخش ۶. عمل کرده (DataFrame با اسم best_pro) با این تفاوت که در اینجا با توجه به مقاله ی salton and buckley مقدار tf و idf را از یک فرمول دیگر به دست می آوریم به این صورت که برای tf از فرمول :

$$0.5 + (0.5 * \frac{tf}{\max(tf)})$$

استفاده میکنیم

و برای idf از فرمول:

$$\log_{10}(N - \frac{n}{n})$$

استفاده میکنیم.

در این قسمت سوال اول assignment 2 تمام میشود.

در بخش ۸. که مربوط به سوال دوم می باشد فایل query ها را مانند فایل cran باز میکنیم.

در بخش ۹. با استفاده از regular exoression اعداد که همان شماره ی داکيومنت ها میباشد و تگ W. را از لیست query ها حذف میکنیم.

سپس همانند داکيومنت ها با استفاده از تگ I. لیست query ها را split میکنیم تا به هر query به صورت جداگانه به وسیله ی index دسترسی داشته باشیم. در ادامه تا رسیدن به بخش

۱۰- کارهایی مانند حذف stopwords ها و stemming انجام شده است. (توضیحات در قسمت کامنت ها)

در بخش ۱۰- یک DataFrame با نام j_query ساخته شد که ستون ها . شماره ی query ها میباشد و سطرها شماره ی داکيومنت ها.

در بخش ۱۱- برای هر خانه ی جدول که یک query را نشان میدهد برای هر query تک تک کلماتش جدا میشوند و مقدار $idf*tf$ میشود این مقدار را برای تک تک کلمات هر query به دست می آوریم و در آخر با هم جمع میکنیم.

یعنی برای هر query مقدار $tf*idf$ کل کلماتش را در هر داکيومنت حساب میکنیم و با هم جمع میکنیم

مقدار tf و idf در بخش ۶- محاسبه شده است.

در بخش ۱۲- نیز مانند بخش ۱۱- یک DataFrame با نام j_query_best ساخته شده با این تفاوت که برای مقدار $tf*idf$ هر خانه از مقادیر tf و idf بخش ۷- استفاده کردیم.

در بخش ۱۳- فایل relevant باز و چاپ شده است.

این لیست شامل شماره ی داکيومنت هایی است که برای تعداد ۱۰ query که در سوال خواسته شده. (این query ها شماره ی ۱۰ تا ۱۹ لیست query ها میباشد)

در بخش ۱۴- میخواهیم مقدار precision و recall را برای هر query با هر دو روش به دست آوریم (روش معمول و یکی از روش های گفته شده در مقاله ی salton & backley)

برای این کار ابتدا برای هر query از ۱۰ query جدولی که مجموع مقدار $tf*idf$ را حساب کردیم (j_query) به صورت نزولی sort میکنیم و با relevant مقایسه میکنیم و شماره هایی از داکيومنت ها که در لیست sort شده آمده و در relevant موجود میباشد را حساب میکنیم و precision و recall را محاسبه و چاپ کرده. (این کار ابتدا برای ۱۰ داکيومنت اول سپس برای ۵۰ و ۱۰۰ و ۵۰۰ داکيومنت sort شده ی اول انجام شد)

این کار را برای روش دوم نیز انجام داده و precisin و recall را محاسبه میکنیم. با مقایسه ی precision در این دو روش مشاهده میشود که روش گفته شده در مقاله که در این کد با نام j_query_best تعریف شده است بهتر عمل میکند. بهتر عمل کرده (یعنی تعداد بیشتری از داکيومنت های لیست relevant را برای هر query برگردانده است).