

بسمه تعالی
گزارش کار پروژه پایانی درس بازیابی اطلاعات
استاد : دکتر حسین امیر خانی
روح الله مظفری
۹۵۱۳۲۰۰۰۵۵

مقدمه :

در این گزارش شیوه جمع آوری اطلاعات از سایت دنیای اقتصاد و همچنین روش tf-idf مورد بررسی قرار خواهد گرفت.

سرفصل ها :

۱- جمع آوری اطلاعات (Crawling)

2-مقایسه داکيومنت های مشابه با روش tf-idf

۱ - جمع آوری اطلاعات (Crawling)

برای جمع آوری دیتا نیاز به دانلود صفحات مختلفی از سایت مورد نظر به شکل html داریم. در این جا از کتابخانه requests زبان پایتون استفاده می کنیم. فرمت url سایت برای صفحات خبری به شکل زیر است

```
https://donya-e-qtasad.com/%D8%A8%D8%AE%D8%B4-%D8%AE%D8%A8%D8%B1-64?np={news_number}
```

که new_number شماره خبر مورد نظر است. تنها نیاز است که این لینک را با کتابخانه requests دانلود کنیم. پس از آن یک

ابجکت از کلاس BeautifulSoup

می سازیم . BeautifulSoup به ما کمک میکند که صفحه دلود شده با به متغییر های قابل دسترسی (tag های html) تغییر شکل دهیم. با کمک متد find_all تمام لینک های قرار گرفته در تگ های h2 که عناوین خبری ما هستند را استخراج میکنیم. این لینک ها ما را به سمت صفحه ی اصلی هر خبر هدایت می کند. برای دالود هر صفحه خبر و استخراج داده از آن از کتاب خانه newspaper استفاده می شود. این کتاب خانه متن و عنوان خبر را پس از دالود کردن صفحه خبر ذخیره میکند. این عملیات را تا پایان پیدا کردن آخرین خبر ادامه می دهیم. دیتای استخراج شده شامل لینک صفحه خبر و عنوان و متن آن می باشد. این دیتا را توسط کتاب خانه pandas به شکل یک ابجت DataFrame تغییر میدهیم. پس از آن این دیتا را در داخل یک فایل با فرمت csv ذخیره میکنم. این قطعه کد در داخل فایل scraper.py قرار دارد. پس از اجرای ان فایل دیتای جمع آوری شده در فایل csv با نام new_data.csv قابل مشاهده خواهد بود.

2-مقایسه داکيومنت های مشابه با روش tf-idf

برای مقایسه و پیدا کردن داکيومنت های مرتبط با کویری کاربر از روش rf-idf استفاده میشود. ابتدا توسط کتاب خانه pandas فایل csv مورد نظر را read میکنیم و آن را به فرمت ارایه تغییر می دهیم. ابتدا یک شی از کلاس TfidfVectorizer می سازیم و تمام داکيومنت های مورد نظر را fit میکنیم. برای تمام کویری های کاربر متد find_query صدا زده میشود. در این جا ابتدا کویری توسط vectorizer زا به شکل tfidf در می اوریم. پس این مرحله ما vector های تمام داکيومنت و کویری کاربر را در اختیار داریم کافیسست که این دو را مقایسه کنیم. به دلیل محدود بودن جواب برگشتی که در این تست عدد ۱۰ ثبت شده این فرایند مقایسه به صورت خطی انجام میشود. مقایسه vector ها توسط متد cosine_similarity انجام میشود. پس از مرتب کردن cosine های بدست آمده ۱۰ داکيومنت با بیشترین مقدار شباهت را باز میگردانیم. در دایرکتوری samples تعدادی کویری همراه با داکيومنت ها و میزان شباهت آنها به ترتیب نزدیک ترین داکيومنت به کویری ذخیره شده است.

کد های مربوطه در صفحه ی گیت هاب به آدرس

<https://github.com/ruhollahmozaferi/tf-idf-method-for-information-retrieval>

قابل مشاهده می باشد.