

# تمرین برنامه‌نویسی شماره ۲

## درس مبانی بازیابی اطلاعات و جستجوی وب

### شاخص‌گذاری معکوس

مهلت ارسال پاسخ: ۱۳۹۶/۰۹/۲۹

#### ۱ مجموعه دادگان<sup>۱</sup>

دادگان در نظر گرفته شده برای این تمرین، مجموعه دادگان متنی Reuters-21578 می‌باشد. این مجموعه دادگان شامل متن تعداد ۲۱۵۷۸ مورد از خبرهایست که در سال ۱۹۸۷ در مجله رویترز به زبان انگلیسی منتشر شده است. این مجموعه دادگان در فرمت SGML (فرمتی مشابه XML دارد) ارائه شده است، شامل تگهای متفاوت با معانی مختلف که بخشهای اطلاعاتی متفاوت را در فایل دادگان از هم مجزا نموده‌اند. جزئیات مربوط به این مجموعه دادگان و قالب و معنی تگهای مختلف آن، در فایل **readme.txt** به همراه فایل‌های اصلی مجموعه دادگان قابل دسترس و مشاهده است. در اینجا برخی از تگها و خصائص مهم که به منظور حل تکلیف جاری و تکالیف بعدی لازم است، شرح داده می‌شود و در صورت نیاز به اطلاعات بیشتر می‌توانید جزئیات کاملتر را از فایل **readme** استخراج نمایید.

مجموعه دادگان Reuters-21578 متشکل از ۲۲ فایل با پسوند **sgm** می‌باشد که هر فایل شامل حدود ۱۰۰۰ قطعه متن متناظر با خبرهای جمع‌آوری شده است. در هر فایل **sgm**، ابتدا و انتهای قطعه مربوط به هر یک از خبرها توسط یک زوج تگ **<REUTERS>....</REUTERS>** شامل تعدادی خصیصه، مشخص شده است. از جمله خصائص مورد نیاز ما در این تمرین، **NEWID** و **TOPICS** است که به ترتیب بیانگر **id** سند (خبر) و اینکه آیا این خبر مرتبط با موضوع خاصی هست (**YES**) یا خیر (**NO**)، می‌باشد. لازم به ذکر است که خبرها ممکن است به یک یا چند موضوع (**topic**) مرتبط باشند.

تگ بعدی که مورد نیاز است **<TOPICS>.....</TOPICS>** است که بیانگر موضوع خبر است. هر کدام از موضوعاتی که خبر با آن مرتبط بوده است توسط محتوای یک زوج تگ **<D>....</D>** به صورت متداخل در این تگ بیان شده است.

---

<sup>1</sup> Dataset

تگ بعدی زوج تگ <TEXT>.....</TEXT> است که شامل اطلاعات متن اصلی خبر می‌باشد. داخل این تگ تعدادی تگ داخلی دیگر وجود دارد از جمله زوج تگ <TITLE>.....</TITLE> که محتوای آن شامل عنوان خبر است همچنین زوج تگ <BODY>.....</BODY> که محتوای آن بیانگر متن اصلی خبر می‌باشد.

## ۲ شرح تمرین

هدف این تمرین آشنایی دانشجویان با مفاهیمی همچون Document Frequency, Term Frequency و شاخص معکوس<sup>۲</sup> در فرایند بازیابی اطلاعات می‌باشد. امروزه کتابخانه‌های آماده‌ای همچون Lucene وجود دارند که بسیاری از محاسبات و عملیات موردنیاز در فرایند بازیابی اطلاعات را انجام داده و خروجی نهایی شامل inverted index را ارائه می‌دهند که قابلیت جستجو و بسیاری از عملیات دیگر را خواهد داشت. از این کتابخانه ها می‌توان به راحتی در پیاده‌سازی یک سیستم بازیابی اطلاعات استفاده نمود و از قابلیت‌های آنها استفاده کرد اما هدف ما ملموس بودن بیشتر این مفاهیم برای دانشجویان می‌باشد. در نتیجه دانشجویان باید فرایند پارس نمودن اولیه محتوای مجموعه دادگان و عملیات پیش‌پردازش متون و تشکیل ساختار inverted index را خودشان به یک زبان برنامه‌نویسی دلخواه، پیاده‌سازی نمایند. مراحل کلی کار به صورت زیر می‌باشد:

۱. قدم اول تجزیه<sup>۳</sup> فایل‌های مجموعه دادگانی است که در اختیار دارید. برای این کار می‌توانید از کتابخانه‌های موجود برای پارس نمودن فایل‌های XML مثل SAX Parser یا Lucene استفاده نمایید و یا اینکه پارسر خاص خودتان را پیاده‌سازی نمایید. لذا کارهایی که باید در مرحله تجزیه فایل‌های خام مجموعه دادگان انجام دهید شامل موارد زیر است:

- تمام فایل‌های SGM را به ترتیب پارس نمایید.
- قطعه اطلاعات متناظر با هر خبر را می‌توان به عنوان یک سند (document) متنی فرض نمود.
- از اسناد (خبرهای) موجود در هر فایل، فقط اسنادی را پردازش کنید که مقدار مشخصه TOPICS در آنها YES باشد. از اسنادی که مرتبط با topic خاصی نیستند صرف نظر نمایید.
- برای هر سند که شرط مورد قبل را داشت، شماره NEWID متناظر با آن و topic‌هایی را که با آنها در ارتباط است، استخراج نمایید.
- عنوان و متن محتوای سند را نیز استخراج نمایید.

۲. پس از استخراج متن خام از فایل‌های دادگان نیاز است که برخی عملیات پیش‌پردازش<sup>۴</sup> روی آنها اعمال نمایید. عملیات پیش‌پردازش روی متون استخراج شده در عنوان و محتوای بدنه خبر اعمال می‌گردد. همانگونه که در مبحث مطرح شده در کلاس دیدیم، عملیات پیش‌پردازی می‌تواند متنوع و به صلاح‌دید شما انجام شود ولی حداقل کارهایی که باید انجام دهید، شامل موارد زیر است:

<sup>۲</sup> Inverted Index

<sup>۳</sup> parse

<sup>۴</sup> Preprocessing

- Tokenizing: استخراج کلمات موجود در متن عنوان و محتوای هر سند و حذف کاراکترهای ناخواسته مثل علائم و ...
- نرمال سازی نمایش تمام کلمات به صورت حروف کوچک انگلیسی
- اعمال عملیات حذف stop word ها با استفاده از لیست stop word های داده شده در فایل stopwords.txt به همراه دیگر فایل های مربوط به این تمرین.
- عملیات stemming: Stemming را می توانید با استفاده از الگوریتمها و پیاده سازیهای موجود از آنها مثل الگوریتم Porter انجام دهید (پیاده سازیهای مختلف از الگوریتم Porter به زبانهای مختلف برنامه نویسی وجود دارد که به راحتی قابل جستجو، دسترسی و بهره برداریست و نیازی نیست خودتان آن را پیاده سازی نمایید).

### ۳. ساخت شاخص معکوس

پس از انجام عملیات پیش پردازش و استخراج token های نهایی، می توان ساختار شاخص معکوس را تشکیل داد. برای انجام این مرحله به ساختار داده ای که قابلیت جستجو در آن و نمایش نتیجه جستجو را داشته باشد، نیاز می باشد. در قالب کلی ساختار داده شاخص معکوس، باید علاوه بر ID اسناد حاوی یک کلمه، به ازای هر سند، تعداد تکرار آن کلمه در سند (TF) نیز محاسبه شده و وجود داشته باشد. همچنین برای هر کلمه موجود در شاخص، تعداد اسناد حاوی آن کلمه (DF) نیز محاسبه و در شاخص معکوس ذخیره شود. ساختار داده ای که در پیاده سازی از آن استفاده می کنید دلخواه است ولی بهتر است از ساختارهای لیست یا map موجود در زبانهای برنامه نویسی باشد که مدیریت آنها ساده تر صورت گیرد.

## ۳ گزارش

۱. معرفی کتابخانه های استفاده شده در بخش های مختلف کار و توضیح مختصری درمورد نحوه استفاده از این کتابخانه ها مثلاً

برای parse کردن فایل های مجموعه دادگان و یا عملیات stemming.

۲. تهیه آمار اولیه در قالب یک جدول از تعداد اسناد مفید استخراج شده از مجموعه دادگان و تعداد کل موضوعات متمایز مرتبط با این اسناد.
۳. تشریح مراحل انجام عملیات پیش پردازش متن و جزئیات مختصری از آن. بدیهیست اگر کار اضافه ای نسبت به حداقل های بیان شده در بخش ۲-۲ از مستند جاری انجام داده اید باید در این بخش شرح دهید تا در ارزیابی مدنظر قرار گیرد. در نهایت توضیح دهید که فرایند پیش پردازشی که انجام داده اید چه تاثیری بر فرایند بازیابی خواهد داشت.

۴. تهیه آماری از تعداد کلمات متمایز موجود در کل مجموعه اسناد پردازش شده در قالب جدولی مقایسه‌ای شامل تعداد کل token قبل و بعد از حذف علائم ناخواسته، تعداد token پس از حذف stop wordها، تعداد token پس از stemming.
۵. تشریح ساختار داده نهایی شاخص معکوس و جزئیات مربوط به آن (نحوه پیاده‌سازی، عملیات قابل انجام روی این ساختار داده، اطلاعات موجود در آن، تعداد کلمات موجود در شاخص و ...)

## ۴ موارد لازم جهت ارزیابی

- کدهای پیاده‌سازی همراه با توضیحات مفید مرتبط با کلاسها و متدها به صورت comment (فقط فایل‌های کد برنامه‌ها را ارسال نمایید تا حجم بسته خیلی زیاد نشود).
- فایل گزارش کار.
- یک فایل متنی مستخرج از شاخص معکوس به این صورت که به ازای هر token موجود در شاخص، یک سطر شامل tokenID, tokenString, DocumentFrequency نوشته شده باشد.
- یک فایل متنی متناظر با شاخص معکوس به این صورت که به ازای هر token موجود در شاخص، به ازای اسنادی که این کلمه در آنها ظاهر شده است سطرهایی داشته و در هر سطر اطلاعاتی به فرمت tokenID, docID, TermFrequency نوشته شده باشد. به عنوان مثال اگر یکی از tokenها کلمه Memphis با ID برابر با ۷۸۶۵ باشد که در اسناد شماره ۲۳، ۴۵۵ و ۱۳۸۵ به ترتیب به تعداد ۵، ۳ و ۱۲.