

# پروژه پایانی درس مبانی بازیابی اطلاعات و جستجوی وب

## استخراج کلمات کلیدی پرتکرار از یک صفحه خبری در Twitter

محمد رضا دلیری (۹۳۱۲۴۳۰۲۳۳)  
مهدی اکبریان رستاقی (۹۳۱۲۴۳۰۴۳۷)

در این پروژه سعی شده تا با استفاده آخرین نسخه کتابخانه Lucene (نسخه 7.2.1) مجموعه‌ای از کلمات کلیدی پرتکرار (داغ) بکار رفته در یک صفحه خبری در شبکه اجتماعی Twitter در و حوزه اقتصادی و اجتماعی استخراج شده و پس از اندیس‌گذاری، مهم‌ترین آن‌ها بازیابی شود. در این بررسی از صفحه رسمی خبرگزاری تسنیم به آدرس @Tasnimnews\_Fa به‌علت متنوع بودن خروجی خبری آن استفاده شده است.

## ۱ اجرای برنامه

اجرای پروژه سه مرحله اصلی دارد، ابتدا لازم است تا توئیت‌ها crawl شده و سپس اندیس‌گذاری و در نهایت بازیابی شود.

### ۱.۱ استخراج توئیت‌ها

عملیات crawling با استفاده از کلاس Crawler انجام می‌شود. نحوه استفاده از این کلاس به‌شرح زیر است:

```
java -cp ProjectNews.jar ir.ac.um.ce.projectnews.crawler.Crawler  
FLAGS
```

پرچم‌های مورد نیاز عبارتند از:

-i: مشخص‌کننده شناسه صفحه‌ی مورد نظر در توئیتر برای استخراج توئیت‌ها

**-s:** مشخص‌کننده تاریخ شروع بازه زمانی مورد نظر برای دریافت توثیت‌ها (با فرمت YYYY-MM-DD)

**-e:** مشخص‌کننده تاریخ پایان بازه زمانی مورد نظر (با فرمت YYYY-MM-DD)

**-m:** محدود کننده حداکثر تعداد پیام های استخراج شده در بازه زمانی مورد نظر

**-p:** مسیر ذخیره فایل CSV محتوی توثیت‌های استخراج شده

**-n:** نام انتخابی برای فایل CSV

به‌عنوان مثال برای استخراج توثیت‌های صفحه‌ی **Tasnimnews\_Fa** در ژوئن ۲۰۱۸ و ذخیره آن در مسیر **out** ورودی باید به‌صورت زیر باشد:

```
java -cp ProjectNews.jar ir.ac.um.ce.projectnews.crawler.Crawler  
-i Tasnimnews_Fa -s 2018-06-01 -e 2018-06-30 -p out/
```

## ۲.۱ اندیس‌گذاری

با توجه به خام بودن داده‌های برنامه (توثیت‌ها) ابتدا باید پیش‌پردازش روی داده‌ها انجام شود. در فرآیند پیش‌پردازش از کلاس **PersianAnalyzer** کتابخانه **Lucene** و لیست **stop words** ارائه شده توسط **آقای وحید خرازی**<sup>۱</sup> استفاده شده است. با استفاده از کلاس **Indexer** می‌توانیم روی داده‌های ورودی فرآیند پیش‌پردازش را انجام داده و سپس آن‌ها را اندیس‌گذاری کنیم. نحوه اجرای این برنامه به‌صورت زیر است:

```
java -cp ProjectNews.jar ir.ac.um.ce.projectnews.search.Indexer  
CORPUS_FILE [INDEX_DIR]
```

**CORPUS\_FILE:** مسیر فایل **corpus** با فرمت **CSV** که شامل توثیت‌ها با فرمت زیر است:

```
ID, Permalink, Text, Date, Retweets, Favorites
```

**INDEX\_DIR:** مسیر مورد نظر برای ذخیره‌سازی **index** ساخته شده توسط برنامه را مشخص می‌کند. این پارامتر اختیاری است و برنامه به‌صورت پیش‌فرض از پوشه **indices** (در کنار فایل اجرایی برنامه) استفاده می‌کند.

<sup>1</sup>Persian (Farsi) Stop Words List

### ۳.۱ دسته‌بندی توئیت‌ها و یافتن کلمات کلیدی پرتکرار

به‌علت مشخص نبودن موضوع توئیت‌های استخراج شده، ابتدا باید توئیت‌ها دسته‌بندی شده و پیام‌های اضافه حذف شوند. برای مشخص کردن دسته‌بندی هر توئیت، از تکنیک جستجو استفاده شده است. بدین صورت که برای هر دسته (اقتصادی و اجتماعی) یک لیست از کلمات کلیدی مرتبط تعریف شده و با جستجوی یک پرس‌وجو شامل این کلمات، پیام‌های مرتبط بر اساس مقدار threshold شناسایی می‌شوند. سپس این پیام‌ها به نمایه‌ی جدیدی منتقل شده و فرآیندهای آتی روی نمایه جدید انجام می‌شود. مقدار threshold از فرمول زیر بدست می‌آید:

$$\frac{minScore + maxScore}{2}$$

لیست کلمات کلیدی مورد استفاده برای دو بخش اقتصادی و اجتماعی در جدول ۱ آمده است. پس از تهیه یک نمایه جدید به‌ازای هر دسته از کلمات که فقط شامل توئیت‌های مرتبط با آن دسته است، فرآیند جستجو برای یافتن پرتکرارترین کلمات (terms) انجام می‌شود. پس از پایان این جستجو، خروجی برنامه به‌صورت یک لیست متنی از کلمات به‌ترتیب تکرار (نزولی) همراه با تعداد تکرار هر کلمه در خروجی استاندارد چاپ می‌شود. این لیست پس از بررسی از نظر ارتباط معنایی با دسته‌بندی، به‌عنوان خروجی نهایی ارائه می‌شود. چگونگی فراخوانی این برنامه به شرح زیر است:

```
java -cp ProjectNews.jar ir.ac.um.ce.projectnews.search.  
Classifier [INDEX_DIR] QUERIES_FILE [RESULTS_COUNT]
```

**INDEX\_DIR:** به محل ذخیره‌سازی index ساخته شده در مرحله قبل اشاره می‌کند. مانند مرحله پیشین، این پارامتر اختیاری است و برنامه به‌صورت پیش‌فرض از پوشه indices (در کنار فایل اجرایی برنامه) استفاده می‌کند.

**توجه:** در صورت تغییر مقدار این پارامتر، لازم است پارامتر سوم (RESULTS\_COUNT) هم مقداردهی شود.

**QUERIES\_FILE:** مسیر یک فایل متنی ساده حاوی کلمات کلیدی مشخص‌کننده را مشخص می‌کند. این فایل یک سطر دارد.

**RESULTS\_COUNT:** تعداد توئیت‌های بررسی شده را مشخص می‌کند. این پارامتر نیز اختیاری بوده و مقدار پیش‌فرض آن، صد هزار است.

**توجه:** در صورت تغییر مقدار این پارامتر، لازم است برای پارامتر اول (INDEX\_DIR) هم مقداری مشخص شود.

جدول ۱: لیست کلمات کلیدی مورد استفاده برای دسته‌بندی توئیت‌ها

دسته‌بندی اقتصادی	<p>سود، ارز، یارانه، واردات، تومان، نزول، سکه، قیمت، بازرگانی، طلا، نرخ، مسکن، رهن، اجاره، مستاجران، برجام، تجاری، صنایع، بازرگانی، بها، کاهش، دلار، افزایش، افزایشی، کلان، اقتصاد، بازار، تک‌نرخ، حقوق، اونس، بهای، سرمایه، تحریم، سرمایه‌گذار، واردکنندگان، خودرو، نرخ، بانک، تجارت، توسعه، سازمان، فروش، هزینه، میانگین، تورم، جواهر، اتحادیه، رانت، اقلام، کالا، میلیون، هزار، رکورد، خودرو، محصولات، نمایندگی، قاچاق، قرارداد، لابی، یارانه، نفت، نفتی، نقدینگی، نقدی، نقد، خانه، رقم، مالیات، مالیاتی، مالک، مالکانی، پول، صادرکنندگان، صادرات، خصوصی، هواپیما، مصرف‌کنندگان، مصرف، سرمایه‌گذاری، گران، گرانی، آماری، آمار، امار، محصول، کاهش، ایرباس، توتال، فساد، بودجه، شاخص، بورس، رشد، تحریم، اوپک، مدیریت، تولید، واردات، وام، بدهکار، نوساز، یورو، پوند، ریال، برنز، نقره، آلومینیوم، آلومینیوم، مس، آهن، آهن، فولاد، بیکار، کارمندان، بهره، قسط، تحریم، میلیارد، اشتغال، تخلف، سفارش، سرمایه، معاملات، معامله، اسکناس، موجودی، مصرف‌کننده، دلال، دلالان، دزد</p>
دسته‌بندی اجتماعی	<p>اجتماعی، فیلم، مردم، دستگیری، مترو، اتهام، فرهنگیان، قتل، آزادی، فیلترینگ، فیلتر، مصرف، پلیس، ایران، دریاچه، ازدواج، طلاق، مهریه، خانواده، خرید، آموزش، پزشکی، پزشک، محرم، مهمانی، ترافیک، روابط، ادارات، اداره، کارگر، مرگ، مترو، سکانس، بهداشت، تحصیل، پلیس، آب، پزشکی، دارو، سلامت، کشور، حج، حاجی، بیمه، دانش‌آموزان، دانش‌آموز، مهمان، تابستان، تعطیلات، محاکمه، قرص، روانگردان، اعتیاد، معتاد، جسد، میوه، بیماری، تهران، مدرسه، آتش‌سوزی، مستمری، سفر، مسافرت، بازنشستگان، فرهنگستان، تصادف، صنایع‌دستی، کارت، ملی، مخدر، نوجوان، استخدام، ایران، معلم، معلمان، صندوق، شهر، استان، زلزله، سیل، جنگل، عید، مسافرت، مسافر، هتل، رستوران، تفریح</p>

## ۲ کامپایل برنامه

سورس برنامه در پوشه `src` قابل مشاهده است. در صورتی تمایل برای کامپایل سورس کد، لازم است از ابزار `Maven` استفاده کرده و فایل `pom.xml` به عنوان فایل تنظیمات به این ابزار داده شود. لازم به ذکر است که دریافت توئیت ها از کتابخانه `GetOldTweets-java` (نسخه 1.2.0) و برای ذخیره آن ها از استاندارد `CSV` با کمک کتابخانه `OpenCSV` استفاده شده است.

## ۳ سورس کد برنامه

این برنامه شامل سه کلاس اصلی زیر می باشد:

**کلاس `Crawler`:** این کلاس همان طور که بیان شد برای دریافت و استخراج توئیت ها تهیه شده است. با توجه به محدودیت های `API` رسمی `Twitter`، کتابخانه `GetOldTweets-java` مورد استفاده قرار گرفته است.

**کلاس `Indexer` و `Writer`:** این دو کلاس در مجموع برای پیش پردازش، ساخت `index` و ذخیره سازی آن در محل مورد نظر کاربر استفاده می شود..

**کلاس `Classifier` و `Searcher`:** وظیفه ی این کلاس ها در کل جداسازی توئیت های مرتبط با هر دسته بندی و سپس مشخص کردن کلمات کلیدی پرتکرار آن دسته بندی می باشد. برای تشخیص توئیت های مرتبط با هر دسته بندی و لیست کلمات آن، از الگوریتم `Okapi BM25` استفاده شده است.

## ۴ نتایج بدست آمده

پس از اجرای فرآیند گفته شده روی صفحه خبرگزاری تسنیم در بازه زمانی سه ماه (ماه های آوریل، مه و ژوئن ۲۰۱۸) برای دو دسته اقتصادی و اجتماعی، کلمات کلیدی پرتکرار در هر ماه مشخص شد که لیست تفکیکی و مرتب شده آن ها (به صورت نزولی) در جداول ۲ و ۳ ارائه شده است.

جدول ۲: پرتکرارترین کلمات کلیدی در خبرهای اقتصادی

رتبه	آوریل ۲۰۱۸	مه ۲۰۱۸	ژوئن ۲۰۱۸
۱	تومان	ایران	بازار
۲	قیمت	تحریم	ارز
۳	دلار	کشور	ایران
۴	کشور	نفت	سرمایه
۵	ارز	برجام	قیمت
۶	افزایش	خصوصی	دلار
۷	بازار	واگن	افزایش
۸	مصرف	اروپایی	تحریم
۹	تحریم	امریکا	نرخ
۱۰	خودرو	بانکی	مصرف
۱۱	ریال	تولید	نفت
۱۲	سکه	روحانی	اروپا
۱۳	فروش	سرمایه	امریکا
۱۴	نرخ	سوئیس	اهن
۱۵	۴۲۰۰	شرکت	دولت

جدول ۳: پرتکرارترین کلمات کلیدی در خبرهای اجتماعی

رتبه	آوریل ۲۰۱۸	مه ۲۰۱۸	ژوئن ۲۰۱۸
۱	اب	دانش آموزان	آموزش
۲	مترو	زندان	مدرسه
۳	مصرف	اتهام	پخش
۴	ایران	ازاد	دانش آموزان
۵	تجربش	آموزش	ایران
۶	تهران	تهران	فوتبال
۷	دریاچه	حمید صفت	پلیس
۸	عذرخواهی	قتل عمد	اذیت
۹	فروشی	ایران	تجاوز
۱۰	پلیس	بازداشت	مدارس