

## گزارش فاز اول پروژه بازیابی اطلاعات

روژینا کاشفی ۹۸۳۱۱۱۸ - لیلی برکتین ۹۸۳۱۰۷۴

### سوال اول

در گام پیش پردازش از سه عملیات استفاده شده است:

#### 1- Normalization

در نرمال سازی واژه‌ها به شکل استاندارد تبدیل می‌شوند برای مثال تشدید و تنوین حذف می‌شود و حروف عربی به فارسی تبدیل می‌شود و نیم فاصله اصلاح می‌شوند.

#### 2- Tokenization

در این مرحله واژه‌ها را به وسیله فاصله از هم جدا می‌کنیم.

#### 3- Stemming

در این مرحله ریشه کلمات را به جای شکل‌های مختلف کلمات نگه‌داری می‌کنیم.

#### 4- Removing Stop words

در مرحله حذف ایست واژه‌ها، واژه‌های تکرار که عموماً معنای خاصی ندارند حذف می‌شود. برای حذف ایست واژه‌ها از لیست ایست واژه‌های کتابخانه هضم استفاده شده است.

به عنوان مثال متن خبر اول به صورت :

به گزارش خبرگزاری فارس، کنفدراسیون فوتبال آسیا (AFC) در نامه ای رسمی به فدراسیون فوتبال ایران و باشگاه گیتی پسند زمان قرعه کشی جام باشگاه‌های فوتبال آسیا را رسماً اعلام کرد. بر این اساس 25 فروردین ماه 1401 مراسم قرعه کشی جام باشگاه‌های فوتبال آسیا در مالزی برگزار می‌شود. باشگاه گیتی پسند بعنوان قهرمان فوتبال ایران در سال 1400 به این مسابقات راه پیدا کرده است. پیش از این گیتی پسند تجربه 3 دوره حضور در جام باشگاه‌های فوتبال آسیا را داشته که هر سه دوره به فینال مسابقات راه پیدا کرده و یک عنوان قهرمانی و دو مقام دومی بدست آورده است. انتهای پیام/

#### • Normalization

تغییرات پس از normalize کردن را با رنگ زرد مشاهده می‌کنیم. در نرمالایز کردن نیم فاصله‌ها درست شده است. نقطه انتها هر کلمه را از کلمه نهایی جدا شده است.

به گزارش خبرگزاری فارس، کنفدراسیون فوتبال آسیا (AFC) در نامه ای رسمی به فدراسیون فوتبال ایران و باشگاه گیتی پسند زمان قرعه کشی جام باشگاه‌های فوتبال آسیا را رسماً اعلام کرد. بر این اساس 25 فروردین ماه 1401 مراسم قرعه کشی جام باشگاه‌های فوتبال آسیا در مالزی برگزار می‌شود. باشگاه گیتی پسند بعنوان قهرمان فوتبال ایران در سال 1400 به این مسابقات راه پیدا کرده است. پیش از این گیتی پسند تجربه 3 دوره حضور در جام باشگاه‌های فوتبال آسیا را داشته که هر سه دوره به فینال مسابقات راه پیدا کرده و یک عنوان قهرمانی و دو مقام دومی بدست آورده است. انتهای پیام/

#### • TOKENIZATION

در این مرحله داکيومنت کلمه کلمه می‌شود و سپس می‌توانیم بررسی کنیم آیا هر کلمه stopword هست یا نه و کلمات را به ریشه خود برگردانیم.

'به', 'گزارش', 'خبرگزاری', 'فارس', 'کنفدراسیون', 'فوتبال', 'آسیا', 'AFC', 'در', 'نامه\200cای', 'رسمی',  
 'به', 'فدراسیون', 'فوتبال', 'ایران', 'و', 'باشگاه', 'گیتی', 'پسند', 'زمان', 'قرعه', 'کشی', 'جام', 'باشگاه\200cهای',  
 'فوتسال', 'آسیا', 'را', 'رسم', 'اعلام', 'کرد', 'بر', 'این', 'اساس', '25', 'فروردین', 'ماه', '1401', 'مراسم', 'قرعه',  
 'کشی', 'جام', 'باشگاه\200cهای', 'فوتسال', 'آسیا', 'در', 'مالزی', 'برگزار', 'می\200cشود', 'باشگاه', 'گیتی',  
 'پسند', 'بعنوان', 'قهرمان', 'فوتسال', 'ایران', 'در', 'سال', '1400', 'به', 'این', 'مسابقات', 'راه', 'پیدا', 'کرده', 'است',  
 'پیش', 'از', 'این', 'گیتی', 'پسند', 'تجربه', '3', 'دوره', 'حضور', 'در', 'جام', 'باشگاه\200cهای', 'فوتسال', 'آسیا', 'را',  
 'داشته', 'که', 'هر', 'سه', 'دوره', 'به', 'فینال', 'مسابقات', 'راه', 'پیدا', 'کرده', 'و', 'یک', 'عنوان', 'قهرمانی', 'و', 'دو',  
 'مقام', 'دومی', 'بدست', 'آورده', 'است', 'انتهای', 'پیام', ]

## • STOP WORDS & STEMMING

مشاهده می‌کنیم که ویرگول یک stop word است که حذف شده است. یا کلمه هر نیز یک کلمه پرتکرار است که حذف شده است.

در مرحله ریشه‌یابی افعال به صورت ریشه نوشته می‌شود.

'گزارش', 'خبرگزاری', 'فارس', 'کنفدراسیون', 'فوتبال', 'آسیا', 'AFC', 'نامه', 'رسمی', 'فدراسیون', 'فوتبال', 'ایران',  
 'باشگاه', 'گیتی', 'پسند', 'زمان', 'قرعه', 'کشید&کش', 'جام', 'باشگاه', 'فوتسال', 'آسیا', 'رسم', 'اعلام',  
 'اساس', '25', 'فروردین', 'ماه', '1401', 'مراسم', 'قرعه', 'کشید&کش', 'جام', 'باشگاه', 'فوتسال', 'آسیا',  
 'مالزی', 'برگزار',  
 'شد&شو', 'باشگاه', 'گیتی', 'پسند', 'بعنوان', 'قهرمان', 'فوتسال', 'ایران', 'سال', '1400', 'مسابقات', 'اس',  
 'گیتی',  
 'پسند', 'تجربه', '3', 'دوره', 'حضور', 'جام', 'باشگاه', 'فوتسال', 'آسیا', 'دوره', 'فینال', 'مسابقات', 'عنوان', 'قهرمانی',  
 'مقام', 'دومی', 'بدست', 'آورده', 'اس', 'انتهای', 'پیام', ]

## سوال دوم

تابع نارنجی از رابطه zipfs و تابع آبی مقادیری که از شاخص بدست آمده را نشان می‌دهد.

هدف این تابع آن است که نشان دهد که هرچقدر فرکانس کمتر شود رنگ آن بیشتر می‌شود و هرچقدر فرکانس بیشتر شود رنگ آن کمتر می‌شود.

ایست واژه‌ها کلمات پرتکرار هستند. پرتکرارترین کلمه در شاخص قبل از حذف ایست واژه‌ها از پرتکرارترین کلمه در شاخص بعد از حذف ایست واژه‌ها، تکرار بیشتری دارد پس همانطور که مشاهده می‌شود قانون zipf's قبل از حذف ایست واژه‌ها عملکرد دقیق‌تری دارد.

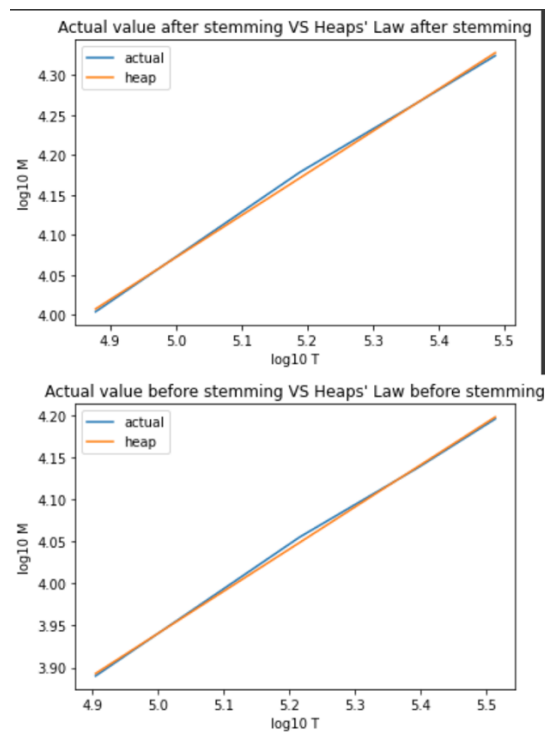


## سوال سوم

ابتدا اندازه vocabulary و collection را برای تعداد مشخصی سند بدست می‌آوریم.

```
After stemming:
First 500 Docs
vocab size : 7767    collection size: 80287
First 1000 Docs
vocab size : 11351   collection size: 164475
First 1500 Docs
vocab size : 13705   collection size: 248561
First 2000 Docs
vocab size : 15679   collection size: 326937
-----
Before stemming
First 500 Docs
vocab size : 10093   collection size: 75476
First 1000 Docs
vocab size : 15082   collection size: 154565
First 1500 Docs
vocab size : 18421   collection size: 233325
First 2000 Docs
vocab size : 21099   collection size: 306798
```

همانطور که مشاهده می‌شود مقادیر بدست آمده از شاخص و مقادیر رابطه heaps تقریباً روی هم افتاده و قانون heaps برای هر دو حالت عملکرد خوبی دارد ولی بعد از ریشه‌یابی عملکرد دقیق‌تری دارد.



## سوال چهارم

ابتدا از کتابخانه هضم استفاده کرده بودم که برای ریشه‌یابی مشکلات زیادی داشت؛ برای مثال بعضی کلمات مفرد را جمع تشخیص میداد (تهران را به تهر تبدیل می‌کرد).

بعضی کلمات که اسم هستند و در آخر آنها حرف م یا ی دارند را فعلی با ضمیر تلقی میکرد (سوم را به سو تبدیل می‌کرد

کتابخانه پارسیوار برای ریشه‌یابی عملکرد بهتری دارد اما نسبت به هضم کندتر است .

فعل مُردَم و اسمِ مُردَم هر دو به مردم تبدیل می‌شود .

بعضی کلمات جمع مثل انتخابات مفرد نمی‌شود و همان انتخابات برگردانده می‌شود .

چند مثال دیگر از اشتباهات ریشه‌یاب پارسیوار:

سال‌ها همان سال‌ها می‌ماند

فعل کردی ریشه‌یابی نمی‌شود.

## سوال پنجم

### الف) پرسمان ساده

کوئری: تحریم‌های امریکا علیه ایران

```
Rank 1:
title: خبرگزاری فارس ۱۹ ساله شد
url: https://www.farsnews.ir/news/14001122000809/خبرگزاری-فارس-۱۹-ساله-شد
-----
Rank 2:
title: سیون فوتبال جمهوری اسلامی ایران هستیم نه جزیره مستقل/ با گفتار ساختارگشانه فدراسیون را به ناکجا آباد می‌برند
url: https://www.farsnews.ir/news/14001117000518/اصول-فدراسیون-فوتبال-جمهوری-اسلامی-ایران-هستیم-نه-جزیره-مستقل-با
-----
Rank 3:
title: احتمال مبادله نازنین زاغری در ازای 530 میلیون دلار
url: https://www.farsnews.ir/news/14001223001080/احتمال-مبادله-نازنین-زاغری-در-ازای-530-میلیون-دلار
-----
Rank 4:
title: متکی: آمریکا با ابزار ناتو به دنبال تجزیه روسیه است
url: https://www.farsnews.ir/news/14001222000749/متکی-آمریکا-با-ابزار-ناتو-به-دنبال-تجزیه-روسیه-است
-----
Rank 5:
title: توفیحات یک منبع آگاه درباره وقفه مذاکرات وین
url: https://www.farsnews.ir/news/14001222000450/توفیحات-یک-منبع-آگاه-درباره-وقفه-مذاکرات-وین
```

باید به گونه‌ای باشد که همه کلمات در وبلاگ‌های مورد نظر باشد و آنکه تعداد تکرار کلمات بیشتری دارد در رتبه بالاتری قرار گیرد. یکی از خروجی‌های را برای صحت عملکرد چک می‌کنیم. بعضی از اخبار مانند خبر زیر کاملاً مرتبط با پرسمان کاربر است ولی در برخی موارد فقط این کلمات را در خبر دارد (ممکن است پراکنده باشد) و از نظر کلیت موضوع مرتبط نیست.

گروه سیاسی خبرگزاری فارس، احد شیرزاد: به نظر می‌رسد مذاکرات وین به روزهای آخر خود نزدیک شده است و براساس اعلام برخی منابع تا رسیدن به توافق راه زیادی باقی نمانده و تنها نیاز است آمریکا به آنچه در برجام برای برداشتن تحریم‌ها وعده داده عمل کند و تضمین معتبری هم برای انجام این تعهدات داده شود.

براساس اعلام مسئولان سیاست خارجه در کنار دریافت تضمین معتبر، برداشتن همه تحریم‌ها از نکات و مطالبات اصلی ایران در مذاکرات است. اما به نظر می‌رسد آمریکا در مقابل عمل به این تعهد برجامی مقاومت می‌کند.

#### \* تحریم‌های آمریکا علیه ایران

آمریکا پس از تسخیر لانه جاسوسی آمریکا در تهران، «کارت‌ر» در تاریخ ۸ نوامبر ۱۹۷۹م. با استناد به قانون «کنترل صدور تسلیحات نظامی»، کشتی حامل لوازم بدکی نظامی متعلق به ایران را توقیف کرد. ارزش این لوازم ۳۰۰ میلیون دلار بود. با اوج‌گیری کشمکش‌های سیاسی بر سر مسئله تصرف سفارت، دولت موقت، اعلام کرد تمام دارایی‌های خود را از بانک‌های آمریکا خارج خواهد کرد. کارت‌ر با لحاظ این احتمال، در کشور شرایط اضطراری اعلام کرد و با استناد به قانون شرایط اضطراری اقتصاد بین‌الملل (IEEPA) و قانون شرایط اضطراری ملی (NEA)، با صدور دستور ویژه‌ی شماره‌ی (۱۲۱۷۰)، تمام دارایی‌های ایران در آمریکا را به تصرف خود درآورد.

### ب) یک پرسمان با عملگر not

کوئری: تحریم‌های امریکا ! ایران

این کوئری کلمات تحریم و امریکا را دارد و ایران را ندارد که تا حد معقولی با کوئری کاربر برابر است.

```

Rank 1:
title: ادامه تحریم های سیاسی علیه المپیک پکن/ژاپن هم به صف منتقدان پیوست
url: https://www.farsnews.ir/news/14001003000306
-----
Rank 2:
title: انتقاد دانشجویان ایرانی در اروپا به برخورد دوگانه مدعیان حقوق بشر با قضایای اوکراین و جنایت های آل سعود
url: https://www.farsnews.ir/news/14001224000014
-----
Rank 3:
title: محور رژیم صهیونیستی از آرمان های نظام اسلامی حذف شده است
url: https://www.farsnews.ir/news/14001222000379
-----
Rank 4:
title: تجربه نشان داده به عهد آمریکا در سده اکر ات نمی شود - اعتماد کرد
url: https://www.farsnews.ir/news/14001203000366
-----
Rank 5:
title: سود مافیای اسلحه سازی آمریکا در دنیا امن بودن جهان است
url: https://www.farsnews.ir/news/14001211000898
-----

```

تصمیم ژاپن برای این موضوع به دنبال **تحریم** دیپلماتیک المپیک زمستانی ۲۰۲۲ پکن از سوی آمریکا، استرالیا، انگلیس و کانادا به خاطر نقض حقوق بشر است. چین از ایالات متحده و سایر کشورها به دلیل نقض بی طرفی سیاسی مورد نیاز در روح منشور المپیک انتقاد کرده است.

## پ) یک پرسمان با عملکرد عبارت کوثری: کنگره ضد تروریست

```

Rank 1:
title: توضیحات یک منبع آگاه درباره وقفه مذاکرات وین
url: https://www.farsnews.ir/news/14001222000450
-----

```

ریگان که نمی خواست کمتر از کنگره **ضد تروریست** جلوه کند، ۳ هفته بعد با صدور دستور ویژه ی (۱۲۶۱۳) ورود هرگونه کالا و خدمات از ایران را ممنوع کرد. وی برای صدور این دستور به بند ۵۰۵ «قانون همکاری های بین المللی امنیتی و توسعه»، مصوب سال ۱۹۸۵م. استناد کرد. در واقع زمانی که ایران درگیر جنگ تحمیلی بود، ایالات متحده به تعریف سیاست های تحریمی پرداخت که بتواند بر نتیجه ی جنگ ایران و عراق تأثیری شگرف بگذارد و رقیب نوپای اسلامی خود را در جنگی نابرابر از میدان به در کند. اما پایان جنگ، تافته ی بافته ی آمریکا را ریش ریش کرد.

فقط یک خبر دقیقاً عبارت کنگره ضد تروریست را دارد. این خبر به طور خاص در مورد کنگره ضد تروریست نیست ولی این عبارت را دارد.

## ت) یک پرسمان پیچیده کوثری: "تحریم هسته ای" آمریکا! ایران

مشاهده می کنیم کوثری فوق شامل phrase و کلمات ساده و عملکرد not است و به همین علت پیچیده است. هیچ خبری نیست که عبارت تحریم هسته ای را داشته باشد ولی ایران را نداشته باشد. برای همین نتایج نامربوط است.

اگر کوثری: "تحریم هسته ای" ایران! آمریکا

اولین خروجی کمی مرتبط است ولی بقیه نامرتب است.

## ث ) یک پرسمان با کلمات نادر

برای پرسمان اورشلیم ! صهیونیست کلمه اورشلیم در شاخص وجود ندارد برای همین پرسمان غزه ! صهیونیست را جستجو میکنیم . غزه نیز کلمه نادر است .

اخبار بدست آمده غزه را دارد ولی صهیونیست را ندارد .

```
Rank 1:
title: سفر تیم فوتبال معلولان فلسطین به تهران
url: https://www.farsnews.ir/news/14001212000749/سفر-تیم-فوتبال-معلولان-فلسطین-به-تهران
-----
Rank 2:
title: بیانیه بسیج ورزشکاران در محکومیت اقدام شرم آور سرمربی تیم ملی امید
url: https://www.farsnews.ir/news/14000929000204/بیانیه-بسیج-ورزشکاران-در-محکومیت-اقدام-شرم-آور-سرمربی-تیم-ملی-امید
-----
Rank 3:
title: اقدام تحسین برانگیز سرمربی الجزایر در حمایت از فلسطین/جام قهرمانی به مردم غزه اهدا شد
url: https://www.farsnews.ir/news/14000928000274/اقدام-تحسین-برانگیز-سرمربی-الجزایر-در-حمایت-از-فلسطین-جام-قهرمانی-به-مردم-غزه-اهدا-شد
-----
Rank 4:
title: اردوغان علاوه بر آرمان فلسطین بر باورهای مردم ترکیه هم پا گذاشت
url: https://www.farsnews.ir/news/14001220000842/اردوغان-علاوه-بر-آرمان-فلسطین-بر-باورهای-مردم-ترکیه-هم-پا-گذاشت
-----
Rank 5:
title: نگاه دوگانه برخی جریان ها به یمن و اوکراین نشانه سرسپردگی شان به آمریکا است
url: https://www.farsnews.ir/news/14001209000540/نگاه-دوگانه-برخی-جریان-ها-به-یمن-و-اوکراین-نشانه-سرسپردگی-شان-به-آمریکا-است
-----
```

تیم ملی فوتبال معلولان قطع عضو فلسطین فوریه گذشته همزمان با روز جهانی معلولان در **غزه** تشکیل شد. این تیم امیدوار است با درخشش در مسابقات غرب آسیا در ترکیه بتواند جواز حضور در مسابقات جام جهانی فوتبال معلولان که قرار است اکتبر آینده در ترکیه برگزار شود، را به دست بیاورد. تیم فوتبال معلولان قطع عضو فلسطین قرار است در گروه غرب آسیا با عراق، ایران، هند و ازبکستان به رقابت بپردازد.