

در گام پیش‌پردازش از سه عملیات استفاده شده است:

### Normalization - 1

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

### Tokenization - 2

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

### Stemming - 3

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

### Removing Stop words - 4

در مرحله حذف ایست واژه‌ها، واژه‌های پر 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 باشگاه های فوتبال آسیا را رسماً اعلام کرد. بر این اساس 25 فروردین ماه 1401 مراسم قرعه کشی جام باشگاه های فوتبال آسیا در مالزی برگزار می شود. 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 باشگاه های 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 نیم‌فاصله‌ها درست شده است. نقطه انتها هر کلمه را از کلمه نهایی جدا شده است.

به گزارش خبرگزاری فارس، کنفدراسیون فوتبال آسیا (AFC) در نامه‌ای رسمی به فدراسیون فوتبال ایران و باشگاه گیتی پسند زمان قرعه کشی جام باشگاه‌های فوتسال آسیا را رسماً اعلام کرد. بر این اساس 25 فروردین ماه 1401 مراسم قرعه کشی جام باشگاه‌های فوتسال آسیا در مالزی برگزار می‌شود. 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 باشگاه‌های 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 /

## • 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\200c\ای', 'رسمی', 'به', 'فدراسیون', 'فوتبال', 'ایران', 'و', 'باشگاه', 'گیتی', 'پسند', 'زمان', 'قرعه', 'کشی', 'جام', 'باشگاه\200c\های', 'فوتسال', 'آسیا', 'را', 'رسماً', 'اعلام', 'کرد', 'بر', 'این', 'اساس', '25', 'فروردین', 'ماه', '1401', 'مراسم', 'قرعه', 'کشی', 'جام', 'باشگاه\200c\های', 'فوتسال', 'آسیا', 'در', 'مالزی', 'برگزار', 'می\200c\شود', '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\200c\های', '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.

'ا 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']

در ادامه به بررسی کوئری هایی برای ارزیابی قابلیت سرچ کوئری می پردازیم:

### الف) پرسمان ساده

کوئری : تحریم‌های آمریکا 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: اصولی: فدراسیون فوتبال جمهوری اسلامی ایران هستیم نه جزیره مستقل/ با گفتار ساختارشکنانه فدراسیون را به ناکجا آباد می‌برد
url: https://www.farsnews.ir/news/14001117000518/اصولی-فدراسیون-فوتبال-جمهوری-اسلامی-ایران-هستیم-نه-جزیره-مستقل-با-گفتار-ساختارشکنانه-فدراسیون-را-به-ناکجا-آباد-می-برد
-----
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: منکی: آمریکا با ابزار ناتو به دنبال تخریب روسیه است
url: https://www.farsnews.ir/news/14001222000749/منکی-آمریکا-با-ابزار-ناتو-به-دنبال-تخریب-روسیه-است
-----
Rank 5:
title: توضیحات یک منبع آگاه درباره وقفه مذاکرات وین
url: https://www.farsnews.ir/news/14001222000450/توضیحات-یک-منبع-آگاه-درباره-وقفه-مذاکرات-وین
-----
```

باید به گونه‌ای باشد که همه کلمات در وبلاگ‌های مورد نظر باشد و آن که 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.

می‌کند. \* **تحریم‌های آمریکا علیه ایران** آمریکا پس از تسخیر لانه جاسوسی آمریکا در تهران، «کارت‌ر» در تاریخ 8 نوامبر 1979م. با استناد به قانون «کنترل صدور تسلیحات نظامی»، کشتی حامل لوازم یدکی نظامی متعلق به ایران را توقیف کرد. ارزش این لوازم 300 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 (12170)، تمام دارایی‌های ایران در آمریکا را به تصرف خود درآورد. مجموع پس‌اندازها و اوراق بهادار بلوکه شده‌ی ایران بالغ بر 12 میلیارد دلار بود. حدود 1.4 میلیارد دلار در بانک «فدرال رزرو» آمریکا، 5.

### ب) یک پرسمان با عملگر NOT

کوئری : تحریم‌های آمریکا ! ایران  
این کوئری کلمات تحریم و آمریکا را دارد و ایران را ندارد که تا حد معقولی با کوئری کاربر برابر است.

```
Rank 1:
title: ادامه تحریم‌های سیاسی علیه المپیک یکن/زاین 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 جنایت‌های آل سعود
url: https://www.farsnews.ir/news/14001224000014
-----
Rank 3:
title: محو رژیم صهیونیستی از آرمان‌های نظام اسلامی حذف نشده است
url: https://www.farsnews.ir/news/14001222000379
-----
Rank 4:
title: تجربه نشان داده به عهد آمریکا در مذاکرات نمی‌شود اعتماد کرد
url: https://www.farsnews.ir/news/14001203000366
-----
Rank 5:
title: سود مافیای اسلحه‌سازی آمریکا در ناامن بودن جهان است
url: https://www.farsnews.ir/news/14001211000898
-----
```

لایحه‌ای به مجلس برده و با تصویب آن واردات نفت از ایران را ممنوع کردند. ریگان که نمی‌خواست کمتر از کنگره **ضدتروریست** جلوه کند، 3 هفته بعد با صدور دستور ویژه‌ی (12613) ورود هرگونه کالا و خدمات از ایران را ممنوع کرد. وی برای صدور این دستور به بند 505 «قانون همکاری‌های بین‌المللی امنیتی و توسعه»، مصوب سال 1985م. استناد کرد. در واقع زمانی که ایران درگیر جنگ تحمیلی بود، ایالات متحده به تعریف سیاست‌های تحریمی پرداخت که بتواند بر نتیجه‌ی جنگ ایران و عراق تأثیری شگرف بگذارد و رقیب نوپای اسلامی خود را در جنگی نابرابر از میدان به در کند. اما پایان جنگ، تافته‌ی بافته‌ی آمریکا را ریش ریش کرد. \* پایان جنگ، تداوم تحریم با پذیرش قطعنامه‌ی (598) توسط

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

تروریست نیست ولی این عبارت را دارد.

در سیستم‌های بازیابی اطلاعات، برای بهبود دقت و کارایی جستجو، از شاخص‌گذاری موقعیتی و وزن‌دهی اسناد استفاده می‌شود. در این بخش، پیاده‌سازی شاخص‌گذاری موقعیتی همراه با محاسبه وزن‌های TF-IDF و ایجاد لیست قهرمانان ارائه شده است. هدف این بخش، پردازش و تحلیل اسناد برای بازیابی سریع و موثر اطلاعات مرتبط با کوئری‌های کاربر است.

محاسبه وزن‌دهی اسناد و عبارات جستجو  
فرآیند وزن‌دهی شامل دو معیار اصلی است:

#### TF (Term Frequency):

نشان‌دهنده تعداد دفعات تکرار یک کلمه در سند است. مقدار TF بر اساس فرمول استاندارد  $\log_{10}(freq) + 1$  محاسبه شده تا تأثیر کلمات پرتکرار کاهش یابد.

#### IDF (Inverse Document Frequency):

میزان اهمیت یک کلمه در کل مجموعه اسناد را تعیین می‌کند. این مقدار طبق رابطه  $\log_{10}(N/n)$  محاسبه می‌شود که در آن N تعداد کل اسناد و n تعداد اسنادی است که شامل آن کلمه هستند.

با استفاده از این دو معیار، مقدار TF-IDF برای هر کلمه در اسناد محاسبه شده و از آن برای وزن‌دهی به عبارات در فرآیند بازیابی استفاده می‌شود.

#### ساخت لیست قهرمانان (Champions List)

برای افزایش کارایی بازیابی اطلاعات، لیستی از مهم‌ترین اسناد برای هر کلمه ساخته می‌شود. این لیست شامل ۲ سند با بالاترین وزن TF-IDF برای هر کلمه است که به کاهش تعداد اسنادی که در فرآیند جستجو بررسی می‌شوند، کمک می‌کند.

#### نرمال‌سازی اسناد

از آنجایی که طول اسناد متفاوت است، برای مقایسه بهتر وزن‌ها، یک بردار نرمال برای هر سند محاسبه می‌شود. در این مرحله، مجموع مربعات وزن‌های کلمات هر سند محاسبه شده و جذر آن به عنوان مقدار نرمال‌سازی‌شده در نظر گرفته می‌شود.

## نتیجه‌گیری

در این بخش، مکانیزم شاخص‌گذاری موقعیتی و وزن‌دهی به اسناد پیاده‌سازی شده است که به بهبود دقت بازیابی و افزایش سرعت پردازش جستجو کمک می‌کند. مراحل انجام‌شده شامل محاسبه وزن‌های TF-IDF، ایجاد لیست قهرمانان و نرمال‌سازی اسناد است که همگی در راستای بهینه‌سازی سیستم جستجو و کاهش هزینه پردازشی انجام گرفته‌اند.

در زیر چند مثال برای بررسی کارایی سیستم جست و جوی خود را بررسی خواهیم کرد.

## ۱- الف) یک کوئری از کلمات ساده و متداول تک کلمه ای:

```
query = 'فوتبال'
query_tokens = preprocess([query], True, True)[0]
scores = cosine_similarity(query_tokens, positional_index, len(contents), docs_norms)
print_result(scores, k=5)
```

✓ 0.0s

1466 title: نکونام: نفتی‌ها بهترین بازی خود را انجام دادند/بازیکنان جدیدمان کیفیت بالای خود را نشان دادند  
url: <https://www.farsnews.ir/news/14001204001165/نکونام-نفتی-ها-بهترین-بازی-خود-را-انجام-دادند-بازیکنان-جدیدمان-کیفیت-بالای-خود-را-نشان-دادند>

81 title: ماجدی: فوتبال کشور به تغییرات نیاز دارد  
url: <https://www.farsnews.ir/news/14001223000539/ماجدی-فوتبال-کشور-به-تغییرات-نیاز-دارد>

6690 title: امیدواری ملی پوش سابق فوتبال ساحلی بابت تغییرات در کادرفنی تیم ملی  
url: <https://www.farsnews.ir/news/14000927000275/امیدواری-ملی-پوش-سابق-فوتبال-ساحلی-بابت-تغییرات-در-کادرفنی-تیم-ملی>

139 title: تقدیر مربی تیم فوتبال خلیج فارس از ماجدی بابت رسیدگی به شائبه تبانی در لیگ جوانان  
url: <https://www.farsnews.ir/news/14001222000297/تقدیر-مربی-تیم-فوتبال-خلیج-فارس-از-ماجدی-بابت-رسیدگی-به-شائبه-تبانی-در-لیگ-جوانان>

860 title: دیدار مدیر تیم ملی ایران با سفیر کره جنوبی/هیان: امیدوارم مقابل ایران مساوی کنیم  
url: <https://www.farsnews.ir/news/14001213000138/دیدار-مدیر-تیم-ملی-ایران-با-سفیر-کره-جنوبی-هیان-امیدوارم-مقابل-ایران-مساوی-کنیم>

همانطور که از تیتراخبار مشخص است همگی مربوط به فوتبال می باشند. به عنوان مثال در تیترا خبر اول کلمه فوتبال مشاهده نمی شود اما اگر خود خبر را بخوانیم بارها کلمه فوتبال را میبینیم:

سرمربی تیم فوتبال فولاد خوزستان درخصوص تغییراتی که در فدراسیون فوتبال ایجاد شده است، تصریح کرد: به عنوان کسی که در فوتبال هستم فقط می‌توانم نظرم را بگویم کاری به تغییرات و اینکه عزیزی خادم چه کاری کرده و چه کاری نکرده ندارم، اما امیدوارم فوتبال ما رو به جلو حرکت کند، چون اگر اینگونه شود برای همه خوب است. نکونام افزود: تجربه نشان داده کسانی که مدیریت فوتبال داشته‌اند و توانسته‌اند کاری بکنند در بدنه فوتبال بوده‌اند. ماجدی آمده و همه ما خوشحال هستیم، چون او هم مثل دادکان فوتبالی است. امیدوارم فوتبال ما رو به جلو حرکت کند. ماجدی را سال‌های زیادی است که می‌شناسیم؛ او رشد کرده و اگر قرار است در فدراسیون بماند که دوست داریم بماند امیدوارم تمام تصمیمات فوتبالی باشد. فدراسیون همه کاره فوتبال است و دوست داریم که میرشاد ماجدی عزیز موفق باشد، چون از بدنه فوتبال است و درد فوتبال را می‌داند. انتهای پیام/.

## ب) یک کوئری از عبارات ساده و متداول چند کلمه ای:

```
query = 'تیم ملی فوتبال'
query_tokens = preprocess([query], True, True)[0]
scores = cosine_similarity(query_tokens, positional_index, len(contents), docs_norms)
print_result(scores, k=5)
✓ 0.0s
```

142 title: اسکوچ: مردم متوجه شده اند که می توانند هدایت تیم ملی را به من واگذار کنند/هاشمیان سواد اروپایی از فوتبال دارد  
url: <https://www.farsnews.ir/news/14001222000329>

2098 title: مصاحبه فارس با کارشناس فوتبال آسیا | از میراث بزرگ کپروش و قدرت ایران با اسکوچ تا انقلاب برانکو در عمان  
url: <https://www.farsnews.ir/news/14001124000522>

1466 title: نگو نام: نفتی ها بهترین بازی خود را انجام دادند/بازیکنان جدیدمان کیفیت بالای خود را نشان دادند  
url: <https://www.farsnews.ir/news/14001204001165>

81 title: ماجدی: فوتبال کشور به تغییرات نیاز دارد  
url: <https://www.farsnews.ir/news/14001223000539>

6690 title: امیدواری ملی پوش سابق فوتبال ساحلی با تغییرات در کادرفنی تیم ملی  
url: <https://www.farsnews.ir/news/14000927000275>

همانطور که از تیتراخبار مشخص است همگی درباره ی تیم ملی فوتبال است.

## پ) یک کوئری دشوار و کم تکرار تک کلمه ای:

```
query = 'واترپلو'
query_tokens = preprocess([query], True, True)[0]
scores = cosine_similarity(query_tokens, positional_index, len(contents), docs_norms)
print_result(scores, k=5)
✓ 0.0s
```

1388 title: تجلیل از خانواده شهید حسن نوفلاح با تقدیم مدال قهرمانی  
url: <https://www.farsnews.ir/news/14001205000938>

5690 title: سرپرست فدراسیون شنا، شیرجه و واترپلو منصوب شد  
url: <https://www.farsnews.ir/news/14001011000202>

5022 title: رضوانی رئیس فدراسیون شنا ماند/ 3 رای سفید برای تنها کاندیدا  
url: <https://www.farsnews.ir/news/14001020000253>

5013 title: برنامه های رئیس فدراسیون شنا برای 4 سال آینده/ رضوانی: تلاش می کنیم به اهدافمان در بازی های آسیایی برسیم  
url: <https://www.farsnews.ir/news/14001020000394>

1300 title: گرفتن ۷۰ نمونه تست دوپینگ در ۸ رشته طی یک هفته  
url: <https://www.farsnews.ir/news/14001207000351>

این کلمه کم تکرار است. همانطور که مشاهده می شود از تیتراخبار مشخص است که مربوط به واترپلو می باشند.

## ت) یک کوئری دشوار و کم تکرار چند کلمه ای:

```
query = 'واکسن کرونا ایرانی'
query_tokens = preprocess([query], True, True)[0]
scores = cosine_similarity(query_tokens, positional_index, len(contents), docs_norms)
print_result(scores, k=5)
✓ 0.0s
```

7937 title: مرندی: رهبر انقلاب دُر سوم واکسن کرونا را دریافت کرده اند  
url: <https://www.farsnews.ir/news/14001117000930>

11966 title: الهیان: وزارت بهداشت درباره علت عدم پیش خرید واکسن فخرآ توضیح دهد  
url: <https://www.farsnews.ir/news/14000729000462>

7228 title: رئیس مسافرت با رعایت اصول بهداشتی بلاشکال است  
url: <https://www.farsnews.ir/news/14001214000475>

9835 title: نایب رئیس مجلس: دولت در تزریق واکسن و کاهش نگرانی های مردم شاهکار کرده است  
url: <https://www.farsnews.ir/news/14000923000300>

9736 title: تأکید مخبر بر حمایت از تولیدکنندگان داخلی واکسن کرونا  
url: <https://www.farsnews.ir/news/14000924000889>



این عبارت کم تکرار است. همانطور که میبینیم نتایج به دست آمده تا حد زیادی مرتبط است. مثلا خبر اول که بالاترین امتیاز را دارد درباره این است که رهبر واکسن ایرانی زده است و خبر دوم درباره ی واکسن فخر است که یک واکسن ایرانی است

## مقایسه با نتایج قبل از پیاده سازی مکانیزم شاخص گذاری موقعیتی و وزن دهی به اسناد:

```
query = 'واکسن کرونا ایرانی'
res = search_query(query)
print_output(res)
```

Rank 1:  
title: شروط حضور هواداران در تمام ورزشگاه های کشور/ مهدی: تغییر زمان دربی نهایی نشده است  
url: <https://www.farsnews.ir/news/14001217000675/>  
-----

Rank 2:  
title: نحوه توزیع سکوی های آزادی برای تماشاگران در دیدار ایران-عراق از زبان خبرنگار عراقی  
url: <https://www.farsnews.ir/news/14001106000376/>  
-----

Rank 3:  
title: وزیر کشور: ممکن است دور سوم واکسن برای سفرهای خارجی اجباری شود  
url: <https://www.farsnews.ir/news/14001217000520/>  
-----

Rank 4:  
title: سخنگوی دولت: ۹۹ درصد سواحلی که در اختیار نهادهای دولتی بود آزاد شد  
url: <https://www.farsnews.ir/news/14001212000835/>  
-----

Rank 5:  
title: همایش علمی دانشگاهیان ایرانی خارج کشور حمایت فارغ التحصیلان ایرانی دانشگاه های برتر دنیا برای بازگشت به کشور  
url: <https://www.farsnews.ir/news/14001212000242/>  
-----

همانطور که مشاهده میکنیم نتایج تقریبا نامربوط اند.

نتایج بدست آمده برای کوئری دشوار و کم تکرار، تقریبا نامرتب است. در ابتدا ما اسنادی را که حاوی کلمات کوئری هستند را بازبایی میکنیم و اسناد را بر اساس اینکه در کدام سند تعداد بیشتری از کلمات کوئری آمده است مرتب میکنیم. برای مثال در کوئری «واکسن کرونا ایرانی» اسنادی را که کلمات واکسن»، «کرونا» و «ایرانی را دارند در رتبه های بالاتر قرار میگیرند و اسنادی که بعضی از این کلمات را ندارند در رتبه های بندی قرار میگیرند.

در پس از اینکه ما پیاده سازی مکانیزم شاخص گذاری موقعیتی و وزن دهی به اسناد انجام دادیم تعداد تکرار یک کلمه در یک سند و اهمیت آن کلمه در کل اسناد را با استفاده از tf-idf پیدا و اسناد می کنیم و با استفاده از شباهت کسینوسی اسناد را از نظر شباهت مقایسه می کنیم. پس هر کلمه یک وزن دارد و اسناد بر اساس وزن و اهمیت کلمات موجود در کوئری بازبایی و رتبه بندی می شوند.

برای همین است که ما در فاز دوم نتایج دقیق تری گرفتیم و توانستیم به خوبی نتایج را رتبه دهی کنیم.