

• مدل سازی اسناد در فضای برداری

برای بازنمایی اسناد به روش بردارهای وزن **tf-idf** کلمات، به دو پارامتر **tf** یا تعداد تکرار هر کلمه در هر سند و **df** یا تعداد اسنادی که هر کلمه در آن ها ظاهر شده است نیاز داریم.

**Tf**:

در هنگام پردازش اسناد در تابع **create\_index**، برای هر سند دیکشنری **termid\_freq={termid:freq,termid2:freq2,...}** که تعداد تکرار کلمات موجود در هر سند را نشان می دهد، بدست می آوریم. سپس هر کدام از این دیکشنری ها را به عنوان **value** در دیکشنری **termid\_freq\_perdoc** قرار می دهیم. در این دیکشنری **key** ها، **docid** ها هستند قرار می دهیم

**termid\_freq\_perdoc = {docid:{termid:freq,termid2:freq2,...},docid2:{termid:freq,termid2:freq2,...}}**

پس بدین طریق تعداد تکرار هر کلمه در هر سند را در اختیار داریم.

**Df**:

در هنگام پردازش اسناد در تابع **create\_index**، در صورت مشاهده شدن هر کلمه در هر سند، به تعداد اسنادی که کلمه در آن ها آمده (**document frequency**) اضافه می کنیم و در نهایت تعداد اسنادی که هر کلمه در آن ها قرار گرفته را در دیکشنری **document\_freq** در اختیار داریم.

**Document\_freq={termid: df,termid2:df2,...}**

بعد از بدست آوردن دو پارامتر **Tf** و **df** برای هر کلمه و سند، دو دیکشنری **termid\_freq\_perdoc** و **document\_freq** را به تابع **tf-idf** برای محاسبه ی وزن هر کلمه در هر سند می دهیم.

در تابع **tf\_idf**، در یک حلقه به ازای هر **docid**، تعداد تکرار هر کلمه ی آن را از دیکشنری **termid\_freq\_perdoc** استخراج می کنیم و در **tf** قرار می دهیم، و تعداد اسنادی که هر کلمه در آن ها ظاهر شده است را نیز از **document\_freq** استخراج می کنیم و در **df** قرار می دهیم و با فرمول زیر وزن هر کلمه را در سند مربوطه بدست می آوریم.

$$w_{t,d} = \log(1 + tf_{t,d}) \times \log_{10}(N / df_t)$$

سپس این وزن ها را در دیکشنری tf\_idf\_weights قرار می دهیم.

$Tf\_idf\_weights = \{docid: \{termid: w, termid2: w2, \dots\}, docid2: \{termid: w, termid2: w2, \dots\}, \dots\}$

برای آنکه از به کار بردن فضای بیش از حد در بازنمایی اسناد جلوگیری شود، تنها کلماتی در دیکشنری termid\_freq\_perdoc تعداد تکرارشان برای هر docid آمده است که در آن docid حداقل یک بار تکرار شده باشند. پس در هنگام محاسبه ی tf-idf چون بر اساس این دیکشنری کلمات موجود برای هر docid وزنشان محاسبه می شود، پس کلماتی که در یک doc وجود ندارند وزنشان محاسبه نمی شود و در بازنمایی اسناد نمی آیند. به این تکنیک index elimination گفته می شود.

خروجی: به دلیل زیاد بودن زمان اجرا برای همه ی خبر ها، 300 خبر از هر فایل خواندم و برنامه را روی آن ها اجرا کردم.

5 کلمه ی با بیش ترین و کم ترین وزن در سند شماره 2:

```
news number 2:
title: اسناد جنگ جهانی دوم در عراق / عراقی نوع مدرسه در جنگ جهانی دوم
5 words with maximum weight in news 2
word: عراقی , weight: 5.18784503214885
word: مدرسه , weight: 4.73518909163542
word: نوع , weight: 4.707562542809043
word: جنگ , weight: 4.453080330679745
word: مدرسه , weight: 4.185899187992144

5 words with minimum weight in news 2
word: کرد , weight: 0.12896523014451372
word: به , weight: 0.4594502051467706
word: کرد , weight: 0.4764004338373566
word: عراقی , weight: 0.49179476132759853
word: جنگ , weight: 0.5243239861545185

Process finished with exit code 0
```

هر چه وزن یک کلمه در سند بیش تر باشد به دو صورت تعبیر می شود:

Tf آن کلمه در آن سند بالا بوده است یعنی کلمه با تعداد تکرار بیش تری در سند ظاهر می شود. هر چه تکرار یک کلمه در یک سند بالاتر باشد وزن آن کلمه در آن سند بالاتر می رود یعنی احتمال این که آن سند به آن کلمه مرتبط باشد بیش تر می شود پس اهمیت کلمه در آن بیش تر است.

Df کلمه کم تر بوده است یا idf آن بیش تر بوده است. یعنی کلمه در اسناد کم تری تکرار شده است که این به معنی خاص تر بودن کلمه است. کلمات رایج که در بیش تر اسناد ظاهر می شوند وزن کلمه در سند را افزایش نمی دهند یعنی احتمال این که این سند به آن کلمه ی خاص مرتبط باشد زیاد نیست. در صورتی که وقتی کلمات خاص در یک سند قرار می گیرند احتمال این که آن سند به آن کلمه مرتبط باشد بیش تر است پس اهمیت کلمه در آن بیش تر است.

- پاسخگویی به پرسمان در فضای برداری

در این بخش ابتدا بردار مخصوص پرسمان ورودی را در تابع  $tf-idf$  محاسبه می کنیم. به این صورت که تمام کلمات پرسمان را استخراج کرده و تعداد تکرار هر کلمه در پرسمان را بدست می آوریم و در دیکشنری  $tf\_query$  قرار می دهیم:  $tf\_query = \{term:tf \text{ in query}, termid2:tf2 \text{ in query}, \dots\}$

و با استفاده از  $df$  آن کلمه که از قبل در دیکشنری  $document\_freq$  داشتیم، وزن کلمات موجود در پرسمان را بدست می آوریم و در دیکشنری  $query\_weights$  قرار می دهیم. سائز بردار پرسمان را نیز محاسبه می کنیم و در  $query\_size$  برای استفاده در محاسبه ی شباهت کسینوسی قرار می دهیم.

$Tf\_query$  شامل  $termid$  تمام کلمات موجود در پرسمان می باشد. پس با یک حلقه در این دیکشنری ، برای تک تک کلمات موجود در پرسمان، ابتدا وزن کلمه در پرسمان را از  $query\_weight$  بدست آورده و در  $w\_t\_q$  قرار می دهیم و سپس  $document$ هایی که کلمه مورد نظر در آن ها حداقل یک بار ظاهر شده است را با استخراج  $postings \text{ list}$  آن ها بدست می آوریم. به این تکنیک  $index \text{ elimination}$  می گوئیم یعنی تنها اسنادی را شباهتشان را با پرسمان محاسبه می کنیم که حداقل یک کلمه ی مشترک با پرسمان داشته باشند و این باعث کم تر شدن محاسبات می شود.

بعد از بدست آوردن این اسناد، وزن کلمه ی مورد نظر در آن ها را با استفاده از دیکشنری  $tf\_idf\_weights$  استخراج می کنیم و در  $w\_t\_d$  ذخیره می کنیم و  $w\_t\_q * w\_t\_d$  را بدست می آوریم و در دیکشنری  $doc\_scores$  ذخیره می کنیم.

$Doc\_scores = \{docid:score, docid2:score2, \dots\}$

این کار را برای تمام کلمات در پرسمان تکرار می کنیم و هر بار برای یک سند مشابه ،  $w\_t\_d * w\_t\_q$  جدید را با  $doc\_scores[docid]$  قبلی جمع می کنیم. به این روش  $term \text{ at a time}$  می گویند و تا آخر حلقه ی روی کلمات پرسمان امتیاز  $doc$ ها به صورت کامل محاسبه نشده است. بعد از تمام شدن این حلقه ی  $for$ ، همه ی امتیاز ها را تقسیم بر  $document\_size[docid] * query\_size$  می کنیم تا شباهت کسینوسی نهایی بدست آید.

سپس از دیکشنری  $doc\_scores$  با استفاده از  $heap$ ، 10 سند با بیش ترین امتیاز را انتخاب کرده و آن ها را به کاربر نمایش می دهیم.

- افزایش سرعت پردازش پرسمان

برای ساخت  $champion \text{ list}$ ، در تابع  $create\_index$  بعد از پردازش اسناد و ساختن شاخص معکوس و دیکشنری  $termed\_freq\_perdoc$  که تعداد تکرار هر کلمه در هر سند را مشخص می کرد، از این دیکشنری

استفاده می کنیم و **champion list** را می سازیم. می دانیم **champion list** لیست اسنادی است که کلمه را با وزن بیش تری در خود دارند. از آن جایی که **idf** یک کلمه ی خاص در تمام اسناد یکسان است پس دلیل تفاوت وزن یک کلمه در اسناد مختلف تعداد تکرار آن کلمه در اسناد یا **tf** آن کلمه در آن سند است. پس به ازای هر کلمه از تمام کلمات (یک **for** در دیکشنری **term\_termid\_match**)، ابتدا روی تمام اسناد (**docid**) موجود در دیکشنری **termed\_freq\_perdoc** یک **for** می زنیم و به ازای هر **docid** دیکشنری تمام کلماتی که در آن سند موجود هستند و تعداد تکرارشان را استخراج می کنیم. سپس چک می کنیم کلمه ی مورد نظر در دیکشنری مرتبط با آن سند است یا خیر، اگر باشد آن سند و تعداد تکرار کلمه در آن را به دیکشنری **champion\_list** در **value** مرتبط با آن کلمه اضافه می کنیم. این دیکشنری بدین صورت است:

$$\text{Champion\_list} = \{\text{termid}:\{\text{docid:tf}, \text{docid2:tf2}, \dots\}, \text{termid2}:\{\text{docid:tf}, \dots\}\}$$

**Docid** های موجود در دیکشنری هر کلمه را براساس تعداد تکرار کلمه در آن ها به ترتیب نزولی مرتب می کنیم. سپس بیش ترین تعداد تکرار هر کلمه در اسنادی کلمه در آن ها ظاهر شده است را به عنوان بیش ترین وزن بدست می آوریم و اسنادی که تعداد تکرار کلمه در آن ها از بیش ترین تکرار-10 کم تر است را از دیکشنری هر **termed** یا به عبارت دیگر از **champion\_list** آن کلمه حذف می کنیم.

سپس در تابع **tf\_idf** برای بدست آوردن اسنادی که می خواهیم شباهت آن ها را با پرسمان محاسبه کنیم به جای **postings list** کلمات در پرسمان از **champion\_list** این کلمات استفاده می کنیم. این کار باعث کاهش محاسبات می شود زیرا تنها شباهت اسنادی را با پرسمان محاسبه می کنیم که احتمال این که در **top k** قرار گیرند به دلیل این که کلمات پرسمان را با وزن بیش تری در خود دارند بیش تر است.