

A hand-drawn diagram of a cell. It features a large, irregular oval shape representing the cell membrane. Inside this oval, there is a smaller, more rounded shape representing the nucleus. To the right of the nucleus, there is a small, dark, hook-like structure labeled with the letter 'C'.

stop word یا کلمات توقیف ، در واقع کلماتی هستند که بطور متداول استفاده می شوند و موتورهای جستجو از جمله گوگل به گونه ای برنامه نویسی شده اند که این کلمات را هم در هنگام ایندکس کردن صفحات وب و هم در هنگام بازیابی آن نادیده بگیرند . از جمله این کلمات می توان به (از، به، و، را، در زبان فارسی اشاره کرد.

در مقام پردازش زبان طبیعی، تمامی این وجود ندارد که در بایگ داده فضا اشتغال
کند. و یا از زمان ارزشمند پردازش غیبه را به خود اختصاص دهد. به همین دلیل، می توان
این کلمات را به راحتی و با ذخیره سازی لیست کلمات توقف حذف کرد (تعداد کم این
کلمات در متن هیچ تأثیری ندارد).

-1

در ابتدا باید گفت آنچه که مورد نیاز یعنی stopword و tokenization را از کتابخانه nltk
ایمپورت کنیم. که stopword برای شناسایی کلمات توقف و web-tokenize برای تبدیل عبارت
کلمات را کفا استفاده می شود. ما باید کتابخانه PorterStemmer را که کارکن تبدیل کلمات
به پایه هست را نیز ایمپورت کنیم. حالا باید رشته ها را وارد کنیم. این که در ابتدا
رشته ها را tokenize و پس از آن stopword را حذف نموده و بعد از آن با stemming
کردن هر رشته، یک لیست کلمات را به کمک termList اضافه می کند. بعد از آنکه
termList برای هر رشته ای ایجاد شد آن را به لیست docs افزوده و پس به سراج رشته بعدی
در docs می رود. نتیجه کفای این عمل یک کتابچه درجده است که docs نام دارد
و شامل termList ها در document است

5

۳- Precision: چند در صد از داده های برگشت آمده با لوری عامه به مقصد رسیدند. (نسبت Relevant Retrieved به Retrieved)

Relevant یعنی از بین تعداد داده‌های موجود، آنهایی که باید برای بیان بیابرد چند درصد را آورده
و چند درصد از Relevant را آورده

در ابتدا متن عادی توکن توکن کرده تا درم ها مشخص شوند.
پس از آن باید جدولی مانند جدول زیر بسازیم

| | T_1 | T_2 |
|-------|-------|-------|
| d_1 | ۱ | ۲ |
| d_2 | ۰ | ۱ |

اکنون نوبت آن رسیده برادر زیر را بسازیم

$$d_1 = (1, 2)$$

$$d_2 = (0, 1)$$

و در نهایت برای مقایسه شباهت ~~کوچک~~ سینوس از فرمول زیر استفاده می کنیم.

$$\text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

Analyzer

از nltk استفاده می کنیم و باید کتابخانه sentiment Intensity را ایمپورت کنیم.
متغیری بنام SIA تعریف و تابع را درون آن قرار می دهیم. حالا اگر ما به SIA یک جمله بدهیم
تا آن جمله را با ما بازیابی کند و یک کلمه کند می توانیم تشخیص دهیم آن جمله منفی است یا مثبت.

negative: میزان منفی بودن که عددی بین ۰ تا ۱ است

positive: میزان مثبت بودن که عددی بین ۰ تا ۱ است

neu: میزان خنثی بودن که عددی بین ۰ تا ۱ است

compound: عددی بین -۱ تا ۱ است در صورتی که منفی باشد جمله منفی و اگر مثبت باشد جمله مثبت است

برای رتبه بندی دایکسونیت ها موجود برای هر کورکی، در نهایت به فرمول TF-IDF رسیدند. یعنی وزنی که قبلاً بدست آوردند (TF) و IDF ضرب شد و فرمول زیر بدست آمد:

$$w_{t,i} = \begin{cases} (1 + \log_{10} F_{t,i}) \times \log_{10} \frac{1}{d_t} & t \in d_i \\ 0 & \text{otherwise} \end{cases}$$

که از این وزن TF-IDF می توان مثلاً فاصله اقلیدس یا کسینوس را بدست آورد تا بفهمیم کدام دایکسونیت ها بر یکدیگر شبیه تر هستند. به کمک ~~TF-IDF~~ TF-IDF scoring می توان امتیاز دایکسونیت را به ازای هر کورکی محاسبه کرد.

$$\text{Score}(q, d_i) = \sum_{t \in q} w_{t,i}$$