

در برنامه مذکور یک کلاس IRSystem تعبیه شده است که در آن document ها تنظیم می شوند و سپس inverted index ها برای کوئری های Boolean و positional index ها برای کوئری های proximity ساخته می شوند. برای تجزیه document نیز، در ابتدا preprocess روی آنها انجام می شود. Preprocess حالت های مختلفی دارد و از لایبری nltk استفاده شده است. در ابتدا ورودی را توکنایز می کند و سپس punctuation ها را از آنها حذف می کند و آنها را تبدیل به حروف کوچک می کند. حال بر حسب پارامتر های تنظیم شده دو کار زیر را انجام می دهد:

- ۱- اگر rm_stop فعال باشد، حروف stop words حذف می شوند مانند to یا is.
 - ۲- اگر do_stem فعال شده باشد، stemming بر روی آنها اتفاق می افتد تا کلمات در یک کلاسه و از یک ریشه به احتمال خوبی با یکدیگر شناخته شوند.
- در بخش build_inverted_index، بر روی document ها حرکت می کنیم و ابتدا بر روی آن document. پیش پردازش گفته شده را انجام می دهیم و سپس به ازای هر term آن، شماره این document را ذخیره می کنیم (posting list)

حال برای هر Boolean query ابتدا بررسی می کنیم که چه نوع سوالی از ما پرسیده شده است:

- ۱- برای مدل and، ابتدا دو ترم اول و سوم را پیش پردازش می کنیم تا مانند ترم های document ها شود، سپس با الگوریتم گفته شده در کلاس intersect sorted، اشتراک دو مجموعه inverted index متناظر با آنها را گرفته و خروجی می دهیم. لازم به ذکر است که برای افزایش سرچ، اولویت را با ارایه ای قرار می دهیم که اندازه اش کم تر است (فرکانس کم تر) چرا که این ارائه محدودیت بیشتری در محاسبه اشتراک خواهد داشت.
- ۲- برای مدل or، مانند and اقدام می کنیم اما به جای اشتراک، اجتماع آنها را می گیریم که چون طبق الگوریتم استفاده شده در ساخت inverted index، می دانیم این دو مجموعه مرتب هستند پس از الگوریتم merge استفاده می کنیم.
- ۳- در حالت not، ترم دوم را پیش پردازش می کنیم و چون inverted index آن sort است، پس به سادگی document هایی که این term را ندارند را محاسبه می کنیم.

نکته حائز اهمیت این است که پارامتر rm_stop و do_stem را چه بگذاریم؟! اگر rm_stop را فعال کنیم در این صورت اگر یکی از term های مورد سوال یک stop words باشد، در این صورت کوئری ارور می گیرد

چرا که پیش پردازش آن خالی می شود. مانند "example or is". پس بهتر است rm_stop فعال نباشد مگر اینکه تضمین شود که ترم های سوال stop words نیستند که با این تضمین طبیعتاً سرعت سرچ بیشتر می شود چون در ساخت inverted index، تعداد ترم کمتری وجود دارد و اعداد به طور کل کوچک تر هستند.

do_stem نیز بهتر است فعال باشد که ترم های از یک ریشه یکی تشخیص داده شوند تا تعداد ترم ها هم کمتر شود و سرعت سرچ بیشتر شود البته که دقت مسئله کاهش می یابد و بستگی به دقت مورد نیاز دارد.

برای سوالات proximity، باید positional index بسازیم چرا که موقعیت ترم ها اهمیت پیدا می کند پس برای هر document پیش پردازش انجام می دهیم و به ازای هر ترم آن، آن document و موقعیت مکانی در آن document را ذخیره می کنیم.

حال به ازای هر کوئری proximity، ابتدا ترم های اول و سوم را برداشت کرده و پیش پردازش می کنیم. از ترم دوم فاصله ماکسیمم قابل قبول را استخراج می کنیم. حال اشتراک document های ترم های اول و دوم را محاسبه می کنیم. (چون positional index هم مرتب است پس از همان الگوریتم intersect sorted استفاده می کنیم.)

حال به ازای هر document مشترک، الگوریتم positional intersect را برای موقعیت های رخداد این دو ترم در این document پیاده می کنیم. اگر موقعیت رخداد این دو کلمه شرط فاصله را رعایت کرده باشد، شماره این document را خروجی می دهیم. در الگوریتم positional intersect، بر روی موقعیت های ترم اول حرکت می کنیم. و به ازای هر موقعیت بررسی می کنیم که موقعیت رخداد ترم دوم نزدیک به آن داریم یا خیر. بدین منظور طبیعتاً موقعیت رخداد ترم دوم که قبلاً بررسی شده است، در قبل از این موقعیت قرار دارد پس اگر فاصله آن بیش تر از kprox (فاصله خواسته شده) موقعیت ترم دوم بعدی را در نظر می گیریم ولی اگر باز هم این فاصله خواسته شده نشود، موقعیت را دوباره بعدی قرار می دهیم. اگر موقعیت ترم دوم از موقعیت فعلی بیشتر شد یا فاصله آن کم تر مساوی با kprox شد، تمام موقعیت دوم هایی که با این موقعیت (ترم اول) شرط مذکور را دارد را ذخیره می کنیم. لازم به ذکر است که این الگوریتم به طور کل بهتر از صرفاً دو حلقه تو در تو بر روی تمام موقعیت های رخداد این دو ترم عمل می کند و از مرتب بودن این دو موقعیت استفاده کردیم.

حال به طور مشابه نکته حائز اهمیت این است که پارامتر rm_stop و do_stem را چه بگذاریم؟! طبیعتاً بهتر است که هیچ کلمه ای در document حذف نشود چرا که فاصله ها را تغییر می دهد و در نتیجه ممکن است نتیجه برای کوئری خاص اشتباه شود. پس بهتر است rm_stop فعال نباشد اما فعال بودن do_stem ایراد

چندانی ندارد. برای مثال با فعال بودن `rm_stop`، نتیجه برای کوئری `"example near/1 test"` اشتباها خروجی دارد چرا که `"to"` حذف می شود و فاصله دو ترم `test` و `example` در `document` سوم یکی کم تر حساب می شود.