

داتاکیومنتیشن فاز دوم و سوم پروژه اول بازیابی اطلاعات

سید امیرمحمد حسینی و فاطمه ده باشی

در فاز دوم و سوم پروژه با توجه به داده‌های استخراج شده در فاز اول و فایل‌های json موجود و با استفاده از tf-idf و مدل Boolean به بازیابی اسناد و رتبه‌بندی آن‌ها می‌پردازیم و با حالات مختلفی این بازیابی را انجام می‌دهیم تا تاثیر آن را در نتیجه بررسی کنیم.

بررسی بازیابی با استفاده از tf-idf

حالت اول: در حالت اول داده‌ها بدون هیچ تغییری و همان متن سند اولیه هستند و همچنین اسنادی که به ما می‌دهد ۲۰ سند اول با توجه به cosine-similarity دو بردار است. نتایج به دست آمده برای متريک‌های ارزیابی به شرح زیر است:

| Query ID | Precision | Recall | P@5 | P@10 | P@15 | NDCG | MAP |
|----------|-----------|--------|-----|------|----------|----------|----------|
| 1 | 0.85 | 0.85 | 1 | 1 | 0.933333 | 0.880159 | 0.827729 |
| 2 | 0.75 | 0.75 | 0.8 | 0.9 | 0.8 | 0.807485 | 0.62035 |
| 3 | 0.75 | 0.75 | 1 | 1 | 0.866667 | 0.893773 | 0.724583 |
| 4 | 0.8 | 0.8 | 1 | 1 | 0.933333 | 0.921731 | 0.778229 |
| 5 | 0.65 | 0.65 | 0.8 | 0.8 | 0.8 | 0.839798 | 0.558152 |
| 6 | 0.9 | 0.9 | 0.8 | 0.9 | 0.933333 | 0.914308 | 0.790145 |
| 7 | 0.9 | 0.9 | 1 | 1 | 1 | 0.982531 | 0.895 |
| 8 | 0.75 | 0.75 | 1 | 1 | 0.8 | 0.902102 | 0.710778 |
| 9 | 0.7 | 0.7 | 0.8 | 0.7 | 0.733333 | 0.734264 | 0.568969 |
| 10 | 0.75 | 0.75 | 1 | 1 | 0.8 | 0.88797 | 0.700458 |
| 11 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.913765 | 0.767531 |
| 12 | 0.9 | 0.9 | 1 | 1 | 1 | 0.962521 | 0.897368 |
| 13 | 0.85 | 0.85 | 1 | 1 | 1 | 0.930268 | 0.844281 |

| 14 | 0.85 | 0.85 | 1 | 1 | 0.933333 | 0.949925 | 0.829966 | | | |
|----|------|----------|-----|-----|-----------|----------|-----------|--|--|--|
| 15 | 0.8 | 0.8 | 0.8 | 0.8 | 0.866667 | 0.945826 | 0.707016 | | | |
| 16 | 0.85 | 0.85 | 0.8 | 0.9 | 0.866667 | 0.914313 | 0.752645 | | | |
| 17 | 0.85 | 0.85 | 1 | 0.9 | 0.933333 | 0.965396 | 0.806038 | | | |
| 18 | 0.6 | 0.6 | 0.6 | 0.7 | 0.733333 | 0.569275 | 0.463507 | | | |
| 19 | 0.55 | 0.55 | 1 | 0.6 | 0.6 | 0.584802 | 0.436679 | | | |
| 20 | 0.85 | 0.85 | 1 | 1 | 0.933333 | 0.907675 | 0.819367 | | | |
| 21 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.927791 | 0.764336 | | | |
| 22 | 0.8 | 0.8 | 0.8 | 0.8 | 0.866667 | 0.896883 | 0.69315 | | | |
| 23 | 0.75 | 0.75 | 1 | 0.9 | 0.8 | 0.916037 | 0.697633 | | | |
| 24 | 0.75 | 0.75 | 1 | 1 | 0.933333 | 0.933932 | 0.736124 | | | |
| 25 | 0.7 | 0.7 | 1 | 0.8 | 0.8 | 0.818745 | 0.614097 | | | |
| 26 | 0.55 | 0.611111 | 0.6 | 0.7 | 0.6 | 0.756277 | 0.460274 | | | |
| 27 | 0.45 | 0.692308 | 0.8 | 0.6 | 0.533333 | 0.875099 | 0.517848 | | | |
| 28 | 0.85 | 0.515152 | 0.8 | 0.8 | 0.8 | 0.474776 | 0.430817 | | | |
| 29 | 0.1 | 0.666667 | 0.2 | 0.1 | 0.133333 | 0.422495 | 0.211111 | | | |
| 30 | 0.45 | 0.428571 | 1 | 0.7 | 0.533333 | 0.434673 | 0.376668 | | | |
| 31 | 0.2 | 0.266667 | 0 | 0 | 0.0666667 | 0.111153 | 0.037924 | | | |
| 32 | 0.25 | 0.294118 | 0.2 | 0.2 | 0.2 | 0.256676 | 0.0894761 | | | |
| 33 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | | |
| 34 | 0.1 | 0.222222 | 0.2 | 0.1 | 0.0666667 | 0.152946 | 0.0487329 | | | |
| 35 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.831771 | 0.990136 | | | |
| 36 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.960412 | 0.978689 | | | |
| 37 | 0.85 | 0.85 | 1 | 1 | 1 | 0.975305 | 0.841796 | | | |

Mean Reciprocal Rank (MRR): 0.9027

ارزیابی کلی:

بطور کلی برای کوئری‌های قبل از کوئری شماره ۲۶ بازیابی اطلاعات خوبی داریم که نشان دهنده بازیابی خوب اطلاعات است، اما از کوئری ۲۶ به بعد که کوئری‌های Boolean هستند در برخی موارد دقت یا فراخوانی پایینی داریم که به دلیل وجود AND و OR و NOT در کوئری است، زیرا تعداد این واژه‌ها تقریباً برابر با تعداد term‌های اصلی کوئری است و به همین دلیل بردار tf-idf کوئری فاصله زیادی از استناد میگیرد و باعث ایجاد خطأ در بازیابی می‌شود.

:Recall و Precision

همانطور که در جدول دیده می‌شود به دلیل اینکه ۲۰ سند برتر توسط tf-idf بازیابی شده است و پاسخ کوئری‌های معیار نیز عموماً ۲۰ سند هستند پس Recall و Precision بجز در چند مورد خاص که تعداد استناد در فایل کوئری‌های معیار با ۲۰ تفاوت دارد با یکدیگر برابر است. اما به طور کلی مقادیر آنها خوب است که نشان دهنده این است که کل استنادی که بازیابی شده به طرز قابل قبولی مرتبط با کوئری هستند.

:P@K

متريک بعدی P@K است که در ۵، ۱۰ و ۱۵ محا سبه شده و در ۲۰ در واقع همان precision می‌شود زیرا کلا ۲۰ سند داریم. با توجه به مقادیر اين معیار در اکثر موارد با افزایش K، دقت کاهش یافته که نشان دهنده اين است که هر چه به اسناد با رتبه پایین‌تر نزدیک می‌شویم احتمال بازگرداندن اسناد غیرمرتبط تو سط سیستم بازیابی بیشتر می‌شود که این نیز با توجه به کم شدن مقدار cosine-similarity در رتبه‌های پایین منطقی است. همچنین از آنجایی که مقدار P@5 که در واقع مهم‌ترین استناد ما هستند (به دلیل اینکه رتبه‌های پایین‌تر از ۵ عموماً مشاهده نمی‌شوند) در اکثر موارد ۱ است به این معنی است که سیستم بازیابی خوب عمل می‌کند.

:MAP

عملکرد سیستم بازیابی با توجه به معیار MAP متوسط است. اعداد به دست آمده در رنج وسیع ۰,۴۳ تا ۰,۸۹ در کوئری‌های قبل از کوئری شماره ۲۶ هستند. این نشان دهنده این است که در برخی کوئری‌ها این سیستم خوب عمل نمی‌کند و در برخی خوب عمل می‌کند با این حال با توجه به میانگین این اعداد که در حدود ۰,۷۵ درصد است می‌توان گفت این سیستم با توجه به معیار MAP عملکرد قابل قبولی دارد.

:MRR

معیار MRR که در آخر با توجه به همه کوئری‌ها محاسبه شده میزان ۰,۹۰۲۷ را دارد که نشان می‌دهد سیستم بازیابی به خوبی و در اولین رتبه‌هایی که بازیابی کرده اسناد مرتبط را برگردانده و این مقدار برای کوئری‌های قبل از کوئری شماره ۲۶، عدد ۱۱ است که نشان می‌دهد تمامی اسنادی که در رتبه ۱ بازگردانده شده اند جزو اسناد مرتبط بوده اند که بسیار خوب است.

:NDCG

در معیار NDCG به هر سند در فایل معیار که سندهای درست هستند یک مقدار relevance داده شده که به این صورت است:

Doc1:20, Doc2:19, ..., Doc20:1

که یک رابطه ساده برای استفاده در NDCG است تا رتبه‌های بالاتر امتیاز بالاتری داشته باشند. این معیار این خوبی را دارد که دیگر تنها مرتبط بودن یا نبودن مهم نیست بلکه میزان آن نیز مهم است. با استفاده از این رابطه و فرمول NDCG اسناد بازیابی شده را با اسناد موجود در فایل کوئری‌های معیار مقایسه کردیم و اعداد به دست آمده در جدول محاسبه شده اند. با توجه به اینکه در اکثر بازیابی‌ها مقدار NDCG بالای 90 درصد است، می‌توان گفت که سیستم بازیابی با دقت خوبی رتبه بندی را انجام داده و علاوه بر مرتبط بودن یا نبودن، میزان این ارتباط را نیز به خوبی متوجه شده است.

در ادامه با اعمال توابع stemming و lemmatization و normalization نتایج را به دست آورده و مقایسه میکنیم:

:Normalization

| Query ID | Precision | Recall | P@5 | P@10 | P@15 | NDCG | MAP |
|----------|-----------|--------|-----|------|----------|----------|----------|
| 1 | 0.85 | 0.85 | 1 | 1 | 0.866667 | 0.87429 | 0.807793 |
| 2 | 0.7 | 0.7 | 0.8 | 0.9 | 0.8 | 0.80143 | 0.578618 |
| 3 | 0.75 | 0.75 | 1 | 1 | 0.866667 | 0.894438 | 0.727679 |
| 4 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.919317 | 0.764436 |
| 5 | 0.65 | 0.65 | 0.8 | 0.8 | 0.8 | 0.839685 | 0.558152 |
| 6 | 0.85 | 0.85 | 0.8 | 0.9 | 0.933333 | 0.911767 | 0.745145 |
| 7 | 0.85 | 0.85 | 1 | 1 | 1 | 0.977511 | 0.85 |
| 8 | 0.75 | 0.75 | 1 | 1 | 0.866667 | 0.902593 | 0.714553 |
| 9 | 0.75 | 0.75 | 0.8 | 0.7 | 0.733333 | 0.770855 | 0.606469 |
| 10 | 0.75 | 0.75 | 1 | 1 | 0.8 | 0.884838 | 0.702847 |
| 11 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.914481 | 0.773951 |
| 12 | 0.9 | 0.9 | 1 | 1 | 1 | 0.962098 | 0.897368 |
| 13 | 0.85 | 0.85 | 1 | 1 | 0.933333 | 0.931115 | 0.841156 |
| 14 | 0.9 | 0.9 | 1 | 1 | 0.866667 | 0.953672 | 0.872049 |
| 15 | 0.8 | 0.8 | 0.8 | 0.8 | 0.866667 | 0.937112 | 0.694442 |
| 16 | 0.85 | 0.85 | 0.8 | 0.9 | 0.866667 | 0.915498 | 0.75767 |
| 17 | 0.85 | 0.85 | 1 | 0.9 | 0.933333 | 0.96463 | 0.806038 |
| 18 | 0.65 | 0.65 | 0.6 | 0.8 | 0.733333 | 0.588594 | 0.498172 |
| 19 | 0.6 | 0.6 | 1 | 0.6 | 0.666667 | 0.629226 | 0.477616 |
| 20 | 0.85 | 0.85 | 1 | 1 | 0.933333 | 0.905072 | 0.816563 |
| 21 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.931476 | 0.764336 |
| 22 | 0.85 | 0.85 | 0.8 | 0.8 | 0.866667 | 0.901524 | 0.737989 |

| | | | | | | | |
|----|------|----------|-----|-----|-----------|----------|-----------|
| 23 | 0.7 | 0.7 | 1 | 0.9 | 0.8 | 0.895631 | 0.661612 |
| 24 | 0.75 | 0.75 | 1 | 1 | 0.933333 | 0.933932 | 0.736124 |
| 25 | 0.7 | 0.7 | 1 | 0.8 | 0.8 | 0.820036 | 0.619397 |
| 26 | 0.55 | 0.611111 | 0.6 | 0.7 | 0.6 | 0.757602 | 0.463021 |
| 27 | 0.45 | 0.692308 | 0.8 | 0.5 | 0.533333 | 0.873532 | 0.513652 |
| 28 | 0.85 | 0.515152 | 0.8 | 0.8 | 0.8 | 0.474271 | 0.430817 |
| 29 | 0.1 | 0.666667 | 0.2 | 0.1 | 0.133333 | 0.380254 | 0.166667 |
| 30 | 0.45 | 0.428571 | 0.8 | 0.7 | 0.533333 | 0.430885 | 0.367156 |
| 31 | 0.2 | 0.266667 | 0 | 0 | 0.0666667 | 0.111153 | 0.037924 |
| 32 | 0.25 | 0.294118 | 0.2 | 0.2 | 0.2 | 0.256676 | 0.0894761 |
| 33 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 34 | 0.1 | 0.222222 | 0.2 | 0.1 | 0.0666667 | 0.152071 | 0.0481481 |
| 35 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.831771 | 0.990136 |
| 36 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.960412 | 0.978689 |
| 37 | 0.85 | 0.85 | 1 | 1 | 1 | 0.975064 | 0.844281 |

Mean Reciprocal Rank (MRR): 0.8982

با توجه به نتایج با اعمال نرمالیزیشن حدود مینیمم و ماکسیمم در معیارها محدودتر شد. در حالتی که اسناد اصلی را داشتیم حدود NDCG بین ۰,۹۶ تا ۰,۹۸ بود اما بعد از نرمالیزیشن حدود آن به ۰,۵۸ تا ۰,۹۷ رسید و همچنین برای MAP این حدود از ۰,۴۳ تا ۰,۸۹ به ۰,۴۶ تا ۰,۸۹ رسید. این به این معنی است که با نرمالیزیشن شاید در کوئری‌هایی که NDCG و MAP بالایی دارند یعنی بسیار خوب عمل کردند دچار مقداری کاهش شویم اما مقدار پیشرفته‌ی که در کوئری‌های با مقادیر NDCG و MAP پایین به دست می‌آید بیشتر است و همچنین در بازیابی‌هایی که دقیق هستند و مقادیر بالا در معیارها دارند این تغییرات جزئی تاثیر زیادی در رتبه‌بندی یا بازیابی موارد مرتبط ندارند اما تغییری که در بازیابی‌های ضعیف‌تر اتفاق می‌افتد بسیار بیشتر است زیرا اکثر اسنادی که داریم ارتباط زیادی با کوئری ندارند و این امر باعث می‌شود هر چه مقدار tf-idf کمتر شود امتیاز اسناد بسیار نزدیک به هم شود و تغییر خیلی کوچک باعث تغییر زیاد در رتبه‌بندی‌ها شود. این مسئله در بازیابی‌های با اسناد بسیار زیاد بیشتر نیز می‌شود. پس به طور کلی می‌توان گفت نرمالیزیشن عملکرد سیستم را بهتر کرد.

:Lemmatization

| Query ID | Precision | Recall | P@5 | P@10 | P@15 | NDCG | MAP |
|----------|-----------|--------|-----|------|----------|----------|----------|
| 1 | 0.55 | 0.55 | 0.8 | 0.8 | 0.666667 | 0.750425 | 0.435857 |
| 2 | 0.3 | 0.3 | 0.8 | 0.6 | 0.4 | 0.43811 | 0.225417 |
| 3 | 0.6 | 0.6 | 0.8 | 0.7 | 0.666667 | 0.638433 | 0.472556 |
| 4 | 0.6 | 0.6 | 1 | 0.8 | 0.666667 | 0.811119 | 0.535571 |
| 5 | 0.7 | 0.7 | 1 | 0.9 | 0.8 | 0.835744 | 0.650229 |
| 6 | 0.8 | 0.8 | 0.8 | 0.9 | 0.866667 | 0.928118 | 0.70661 |
| 7 | 0.8 | 0.8 | 1 | 0.9 | 0.8 | 0.93739 | 0.738336 |
| 8 | 0.7 | 0.7 | 1 | 0.9 | 0.8 | 0.896385 | 0.637387 |
| 9 | 0.55 | 0.55 | 0.8 | 0.6 | 0.533333 | 0.662788 | 0.399556 |
| 10 | 0.65 | 0.65 | 0.8 | 0.9 | 0.733333 | 0.848066 | 0.574169 |
| 11 | 0.75 | 0.75 | 1 | 1 | 0.866667 | 0.888215 | 0.713829 |
| 12 | 0.7 | 0.7 | 1 | 1 | 0.8 | 0.893819 | 0.662068 |
| 13 | 0.85 | 0.85 | 1 | 1 | 1 | 0.928563 | 0.844281 |
| 14 | 0.65 | 0.65 | 1 | 0.8 | 0.666667 | 0.825619 | 0.560226 |
| 15 | 0.4 | 0.4 | 0.8 | 0.6 | 0.466667 | 0.594128 | 0.312532 |
| 16 | 0.85 | 0.85 | 0.8 | 0.9 | 0.866667 | 0.921906 | 0.766483 |
| 17 | 0.65 | 0.65 | 1 | 0.9 | 0.866667 | 0.914037 | 0.623315 |
| 18 | 0.7 | 0.7 | 1 | 0.8 | 0.8 | 0.682114 | 0.630435 |
| 19 | 0.65 | 0.65 | 0.8 | 0.7 | 0.733333 | 0.648645 | 0.50808 |
| 20 | 0.5 | 0.5 | 0.2 | 0.3 | 0.466667 | 0.433444 | 0.235066 |
| 21 | 0.8 | 0.8 | 1 | 1 | 0.866667 | 0.912645 | 0.772843 |
| 22 | 0.7 | 0.7 | 0.8 | 0.8 | 0.866667 | 0.817819 | 0.584427 |
| 23 | 0.7 | 0.7 | 1 | 0.9 | 0.733333 | 0.871646 | 0.6246 |
| 24 | 0.6 | 0.6 | 1 | 0.9 | 0.733333 | 0.764402 | 0.552675 |
| 25 | 0.6 | 0.6 | 0.8 | 0.8 | 0.6 | 0.785801 | 0.478558 |

| | | | | | | | |
|----|------|-----------|-----|-----|-----------|-----------|------------|
| 26 | 0.5 | 0.5555556 | 0.6 | 0.7 | 0.666667 | 0.713917 | 0.427609 |
| 27 | 0.45 | 0.692308 | 0.8 | 0.6 | 0.533333 | 0.879312 | 0.524716 |
| 28 | 0.85 | 0.515152 | 0.8 | 0.8 | 0.8 | 0.463353 | 0.442511 |
| 29 | 0.1 | 0.666667 | 0.2 | 0.1 | 0.133333 | 0.380254 | 0.166667 |
| 30 | 0.45 | 0.428571 | 0.8 | 0.7 | 0.533333 | 0.433967 | 0.367156 |
| 31 | 0.2 | 0.266667 | 0 | 0 | 0.0666667 | 0.110892 | 0.0372222 |
| 32 | 0.05 | 0.0588235 | 0 | 0 | 0.0666667 | 0.0469916 | 0.00392157 |
| 33 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 34 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 35 | 0.7 | 1 | 1 | 1 | 0.8 | 0.846743 | 0.974002 |
| 36 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.963138 | 0.978689 |
| 37 | 0.65 | 0.65 | 1 | 1 | 0.8 | 0.855406 | 0.638235 |

Mean Reciprocal Rank (MRR): 0.8775

همانطور که مشاهده می‌شود تمامی معیارها به طرز قابل توجهی افت پیدا کرد و این به دلیل این است که استاد ما مربوط به مباحث روانشناسی و پزشکی است و این قبیل اسناد دارای تعداد زیادی واژه تخصصی و علمی است و با انجام lemmatization ما این کلمات را عوض می‌کنیم و این باعث می‌شود تاثیر کلمات کلیدی در بازیابی به شدت کاهش یابد و نتایج ضعیفی به دست آید.

Stemming

| Query ID | Precision | Recall | P@5 | P@10 | P@15 | NDCG | MAP |
|----------|-----------|--------|-----|------|-----------|-----------|-----------|
| 1 | 0.4 | 0.4 | 0.4 | 0.5 | 0.533333 | 0.575642 | 0.24895 |
| 2 | 0.3 | 0.3 | 0.6 | 0.5 | 0.333333 | 0.394679 | 0.201131 |
| 3 | 0.35 | 0.35 | 0.6 | 0.5 | 0.4 | 0.548295 | 0.235866 |
| 4 | 0.5 | 0.5 | 1 | 0.8 | 0.6 | 0.759049 | 0.467225 |
| 5 | 0.65 | 0.65 | 1 | 1 | 0.866667 | 0.821658 | 0.65 |
| 6 | 0.7 | 0.7 | 0.8 | 0.9 | 0.866667 | 0.863099 | 0.597422 |
| 7 | 0.5 | 0.5 | 0.8 | 0.5 | 0.533333 | 0.636798 | 0.378875 |
| 8 | 0.15 | 0.15 | 0.2 | 0.1 | 0.0666667 | 0.220758 | 0.0637771 |
| 9 | 0.3 | 0.3 | 0 | 0.4 | 0.333333 | 0.268663 | 0.0952579 |
| 10 | 0.55 | 0.55 | 1 | 0.7 | 0.533333 | 0.712188 | 0.44087 |
| 11 | 0.7 | 0.7 | 1 | 0.9 | 0.733333 | 0.857523 | 0.635237 |
| 12 | 0.15 | 0.15 | 0.4 | 0.3 | 0.2 | 0.303914 | 0.115 |
| 13 | 0.2 | 0.2 | 0 | 0.2 | 0.2 | 0.090766 | 0.0394946 |
| 14 | 0.55 | 0.55 | 1 | 0.8 | 0.666667 | 0.765734 | 0.504371 |
| 15 | 0.1 | 0.1 | 0.2 | 0.2 | 0.133333 | 0.208344 | 0.0642857 |
| 16 | 0.85 | 0.85 | 0.8 | 0.9 | 0.866667 | 0.921966 | 0.766483 |
| 17 | 0.65 | 0.65 | 1 | 0.9 | 0.866667 | 0.91394 | 0.623315 |
| 18 | 0.5 | 0.5 | 0.8 | 0.6 | 0.6 | 0.578512 | 0.362063 |
| 19 | 0.6 | 0.6 | 0.8 | 0.8 | 0.666667 | 0.640839 | 0.493608 |
| 20 | 0.05 | 0.05 | 0.2 | 0.1 | 0.0666667 | 0.0892904 | 0.05 |
| 21 | 0.25 | 0.25 | 0.6 | 0.4 | 0.333333 | 0.533742 | 0.194231 |
| 22 | 0.65 | 0.65 | 0.8 | 0.8 | 0.733333 | 0.754966 | 0.529342 |
| 23 | 0.15 | 0.15 | 0.2 | 0.2 | 0.133333 | 0.11755 | 0.0333333 |
| 24 | 0.1 | 0.1 | 0 | 0.1 | 0.133333 | 0.086366 | 0.0132479 |
| 25 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 | 0.179708 | 0.0657051 |

| | | | | | | | |
|----|------|----------|-----|-----|-----------|-----------|------------|
| 26 | 0.15 | 0.166667 | 0.2 | 0.2 | 0.2 | 0.19005 | 0.0396825 |
| 27 | 0.6 | 0.923077 | 0.8 | 0.8 | 0.733333 | 0.928241 | 0.779446 |
| 28 | 0.95 | 0.575758 | 1 | 0.9 | 0.933333 | 0.539917 | 0.549094 |
| 29 | 0.1 | 0.666667 | 0.2 | 0.1 | 0.133333 | 0.380254 | 0.166667 |
| 30 | 0.5 | 0.47619 | 0.8 | 0.7 | 0.533333 | 0.449732 | 0.390965 |
| 31 | 0.2 | 0.266667 | 0 | 0 | 0.0666667 | 0.110346 | 0.0372222 |
| 32 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 33 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 34 | 0.05 | 0.111111 | 0 | 0 | 0.0666667 | 0.0516568 | 0.00854701 |
| 35 | 0.7 | 1 | 1 | 1 | 0.8 | 0.852466 | 0.974002 |
| 36 | 0.7 | 1 | 1 | 1 | 0.933333 | 0.963138 | 0.978689 |
| 37 | 0.4 | 0.4 | 0.8 | 0.5 | 0.533333 | 0.612673 | 0.326677 |

Mean Reciprocal Rank (MRR): 0.7540

به همان دلیلی که در lemmatization مطرح شد و همچنین افزون بر آن ایجاد کلمات بی معنی و محدود کردن بسیار زبان در stemming که در اسناد مربوط به مباحث علمی نتیجه خوبی را به ارمغان نخواهد داشت، بعد از stemming نیز بسیار بدتر lemmatization نتایج stem کردیم حتی از نتایج lemmatization شد.

بررسی بازیابی با استفاده از Boolean retrieve

نتایج بازیابی با اسناد اولیه و بدون تغییر به شرح زیر است:

| Query ID | Precision | Recall | MAP |
|----------|-----------|--------|----------|
| 26 | 0.352941 | 1 | 0.285643 |
| 27 | 1 | 1 | 1 |
| 28 | 1 | 1 | 1 |
| 29 | 1 | 1 | 1 |
| 30 | 1 | 1 | 1 |
| 31 | 1 | 1 | 1 |
| 32 | 1 | 1 | 1 |
| 33 | 1 | 1 | 1 |
| 34 | 1 | 1 | 1 |
| 35 | 1 | 1 | 1 |
| 36 | 1 | 1 | 1 |

Mean Reciprocal Rank (MRR): 0.9242

در کوئری شماره ۲۶ عملکرد ضعیفی وجود دارد اما با بررسی که انجام شد مشکل از سیستم بازیابی نیست بلکه کوئری های معیار به درستی وارد نشده اند و به همین دلیل این اعداد وجود دارد اما در بقیه کوئری ها همانطور که مشاهده می شود سیستم به خوبی عمل کرده و دقت و فراخوانی ۱ است.

Normalization

| Query ID | Precision | Recall | MAP |
|----------|-----------|----------|----------|
| 26 | 0.352941 | 1 | 0.285643 |
| 27 | 1 | 1 | 1 |
| 28 | 0 | 0 | 0 |
| 29 | 1 | 0.666667 | 0.666667 |
| 30 | 0 | 0 | 0 |
| 31 | 0 | 0 | 0 |
| 32 | 1 | 1 | 1 |
| 33 | 1 | 1 | 1 |
| 34 | 1 | 1 | 1 |
| 35 | 1 | 1 | 1 |
| 36 | 1 | 1 | 1 |

Mean Reciprocal Rank (MRR): 0.6515

بعد از نرمالیزیشن عملکرد بازیابی کاهش پیدا کرده و در بعضی کوئری‌ها هیچ سندی بازیابی نشده است که این به این دلیل است که با در جست و جوی boolean که فقط بودن یا نبودن term‌ها مهم است با تغییر کلمات اسناد ممکن است هیچ سندی بازگردانده نشود یا تعداد اسناد کمتری از تعداد اسناد واقعی مرتبط بازگردانده شود.

Lemmatization

| Query ID | Precision | Recall | MAP |
|----------|-----------|----------|----------|
| 26 | 0.346154 | 1 | 0.27739 |
| 27 | 1 | 1 | 1 |
| 28 | 0 | 0 | 0 |
| 29 | 1 | 0.666667 | 0.666667 |
| 30 | 0 | 0 | 0 |
| 31 | 0 | 0 | 0 |
| 32 | 1 | 0.411765 | 0.411765 |
| 33 | 1 | 1 | 1 |
| 34 | 1 | 0.444444 | 0.444444 |
| 35 | 0.482759 | 1 | 0.653386 |
| 36 | 1 | 1 | 1 |

Mean Reciprocal Rank (MRR): 0.6515

همینطور که مشاهده می‌شود بعد از lemmatization نتایج از normalization هم بدتر می‌شود زیرا کلمات ساده‌تر از حالت قبل می‌شوند و این باعث می‌شود کلمات بیشتری از اسناد ما تغییر کند که در امتداد آن بازیابی مرتبط کمتری اتفاق می‌افتد و در کوئری‌های ۳۲، ۳۴ و ۳۵ از مرحله قبل بدتر می‌شود پس عملکرد سیستم را ضعیف‌تر می‌کند.

• Stemming

| Query ID | Precision | Recall | MAP |
|----------|-----------|----------|----------|
| 26 | 0 | 0 | 0 |
| 27 | 1 | 0.615385 | 0.615385 |
| 28 | 0 | 0 | 0 |
| 29 | 1 | 0.666667 | 0.666667 |
| 30 | 0 | 0 | 0 |
| 31 | 0 | 0 | 0 |
| 32 | 0 | 0 | 0 |
| 33 | 0 | 0 | 0 |
| 34 | 0 | 0 | 0 |
| 35 | 0.482759 | 1 | 0.653386 |
| 36 | 1 | 1 | 1 |

Mean Reciprocal Rank (MRR): 0.3636

بعد از stemming به دلیل اینکه اسناد اولیه علاوه بر ساده شدن و محدود شدن کلمات زبان و احتمال بیشتر عدم تطابق با term های کوئری، کلمات بی معنی نیز ایجاد می شود و این عملکرد را از همه حالات ضعیف تر می کند و همانطور که مشاهده می شود در اکثر کوئری ها عملا هیچ سندی بازیابی نشده است.

بررسی گسترش کوئری با استفاده از تزاروس و تاثیر آن در بازیابی

چهار روش اصلی گسترش کوئری را پیاده سازی کرده ایم:

- ۱- با استفاده از گرفتن متراffد ها از یک تزاروس آماده (ساخته نشده از روی اسناد پروژه).
- ۲- با استفاده از گرفتن متراffد ها از یک تزاروس آماده و همچنین چک کردن وجود داشتن آن در اسناد پروژه.
- ۳- با استفاده از cosine similarity کلمات کوئری و متراffد های آن و گذاشتن threshold.
- ۴- با استفاده از ماتریس هم رخدادی که از روی کلمات اسناد پروژه ساخته شده است.

مثال:

کوئری : Serum HIV testing

کوئری های گسترش داده شده توسط روش اول:

['Serum HIV testing', 'Serum HIV screen', 'Serum HIV examination', 'Serum HIV prove', 'Serum HIV try_out', 'Serum HIV essay', 'Serum HIV try', 'Serum HIV examine', 'Serum HIV quiz', 'Serum HIV test']

کوئری های گسترش داده شده توسط روش دوم:

['serum hiv examination', 'serum hiv test', 'serum hiv testing', 'serum hiv screen']

کوئری های گسترش داده شده توسط روش سوم:

['Serum HIV testing', 'Serum HIV examination', 'Serum HIV test']

کوئری های گسترش داده شده توسط روش چهارم:

['Serum HIV diagnostic', 'Serum HIV criteria']

حالا با محاسبه معیارهای ارزیابی برای کوئری‌های گسترش داده شده کوئری با بیشترین امتیاز را با کوئری اصلی مقایسه می‌کنیم:

| | | | | | | | |
|---|-----|-----|---|---|----------|----------|-----|
| 5 | 0.7 | 0.7 | 1 | 1 | 0.933333 | 0.873713 | 0.7 |
|---|-----|-----|---|---|----------|----------|-----|

سطر بالا مربوط به کوئری ۵ یا همان Serum HIV testing است که در اولین جدول داکیومنت نیز موجود است.

این سطر مربوط به کوئری شماره ۴ از لیست کوئری‌های گسترش داده شده توسط روش اول است که کوئری Serum HIV prove است و همانطور که مشاهده می‌شود در NDCG که قوی ترین معیار برای ارزیابی بازیابی رتبه‌بندی شده است امتیاز بالاتری دارد و نشان می‌دهد که توسط گسترش کوئری توانسته ایم کوئری‌ای پیدا کنیم که از کوئری اصلی بهتر است و می‌توانیم آن را به جای کوئری اصلی استفاده کنیم تا نتایج بهتری داشته باشیم.

جدول زیر کوئری‌های گسترش یافته توسط روش اول را نشان می‌دهد:

| Query ID | Precision | Recall | P@5 | P@10 | P@15 | NDCG | MAP |
|----------|-----------|--------|-----|------|----------|----------|----------|
| 1 | 0.65 | 0.65 | 1 | 1 | 0.8 | 0.841157 | 0.613496 |
| 2 | 0.7 | 0.7 | 1 | 1 | 0.933333 | 0.874084 | 0.7 |
| 3 | 0.7 | 0.7 | 1 | 0.9 | 0.733333 | 0.864632 | 0.622672 |
| 4 | 0.7 | 0.7 | 0.8 | 0.8 | 0.8 | 0.887621 | 0.602453 |
| 5 | 0.7 | 0.7 | 0.8 | 0.9 | 0.8 | 0.862955 | 0.620558 |
| 6 | 0.7 | 0.7 | 1 | 0.9 | 0.8 | 0.863231 | 0.642758 |
| 7 | 0.7 | 0.7 | 1 | 0.8 | 0.666667 | 0.846744 | 0.578995 |
| 8 | 0.7 | 0.7 | 0.8 | 0.8 | 0.8 | 0.844742 | 0.609403 |
| 9 | 0.7 | 0.7 | 1 | 1 | 0.933333 | 0.874084 | 0.7 |
| 10 | 0.7 | 0.7 | 1 | 1 | 0.933333 | 0.874084 | 0.7 |

همانطور که مشاهده می‌شود برخی از این کوئری‌ها NDCG بهتری از کوئری اصلی داشته‌اند و می‌توانیم با استفاده از آنها و یا پیشنهاد آنها به کاربر بازیابی بهتری انجام دهیم و این امر نشان‌دهنده این است که گسترش کوئری می‌تواند تا چه حد در قوی بودن سیستم بازیابی موثر باشد و نتایج را بهتر کند.

لازم به ذکر است که این تنها یک کوئری مثال بود و با بیشتر شدن کوئری‌ها قطعاً نتایج بهتر و بیشتری به دست خواهد آمد که نشان‌دهنده تاثیر گسترش کوئری در خوب بودن سیستم بازیابی است.