

## به نام یکدانه دردانه

سید عارف طباطبایی 9831040

فاز 1:

زیر بخش 1:

ابتدا با `json.load` داک‌های خبری را می‌خوانیم و آن‌ها را به صورت `tuple` های سه تایی شامل `title`، `url` و `content` در می‌آوریم.

در هر `content`، ابتدا `punctuation` ها را حذف می‌کنیم. سپس با استفاده از کتابخانه `parsivar`، متن مربوطه را ابتدا `normalize` و سپس `tokenize` می‌کنیم. سپس توکن‌ها را ریشه‌یابی می‌کنیم و در ادامه نیز `stopwords` را حذف می‌کنیم. برای مثال خواهیم داشت:

به ازای متن ورودی:

'!تقدیم به تو . امیدوارم امروز حالت خوب باشه ؟ سلام خوبی'

نتیجه:

'سلام', 'خوبی', 'امیدوار', 'امروز', 'حالت', 'خوب', 'باشه', 'تقدیم', 'تو'

زیر بخش 2:

ساختمان شاخص مکانی را با استفاده از تابع `defaultdict` در کتابخانه `collections` ایجاد می‌کنیم:

```
# Create the inverted index
index = defaultdict(lambda: {'num': 0, 'positions': defaultdict(lambda: [0, []])})
```

شامل `num` و `positions` که در اینجا `num` به تعداد بار استفاده از یک کلمه در تمام `doc` ها و `positions` به تعداد بار استفاده از آن کلمه در هر `doc` و موقعیت مکانی‌های آن‌ها اشاره دارند.

مطابق تکه کد زیر تمام `doc` ها را زیر و رو کرده و شاخص مکانی را می‌سازیم:

```
for i, doc in preprocessed_data.items():
    if int(i) % 1000 == 0:
        print(f'{int(i)} have been processed.')
    content = doc['content']
    for j, token in enumerate(content):
        index[token]['num'] += 1
        index[token]['positions'][i][0] += 1
        index[token]['positions'][i][1].append(j)
```

در انتها مطابق زیر برای مثال، شاخص مکانی را برای کلمه خبر چاپ می‌کنیم:

```
term = index['خبر']
print(f"Total frequency: {term['num']}")
print("Positions:")
for doc_id, positions in term['positions'].items():
    print(f" Document {doc_id}: total_number: {positions[0]}, positions: {positions[1]}")
print()
```

نتیجه:

```
Total frequency: 1881
Positions:
Document 14: total_number: 1, positions: [6]
Document 19: total_number: 1, positions: [13]
Document 29: total_number: 1, positions: [420]
Document 48: total_number: 1, positions: [360]
Document 56: total_number: 1, positions: [131]
Document 78: total_number: 2, positions: [12, 25]
Document 79: total_number: 1, positions: [55]
Document 135: total_number: 1, positions: [26]
Document 142: total_number: 1, positions: [838]
Document 148: total_number: 2, positions: [320, 682]
Document 150: total_number: 1, positions: [22]
Document 200: total_number: 1, positions: [31]
Document 203: total_number: 1, positions: [67]
Document 206: total_number: 1, positions: [26]
Document 211: total_number: 2, positions: [18, 61]
Document 216: total_number: 2, positions: [124, 135]
Document 218: total_number: 2, positions: [32, 93]
Document 219: total_number: 1, positions: [80]
Document 220: total_number: 1, positions: [94]
Document 222: total_number: 1, positions: [25]
Document 223: total_number: 2, positions: [62, 86]
Document 226: total_number: 1, positions: [13]
```

```
Document 236: total_number: 1, positions: [122]
Document 238: total_number: 1, positions: [154]
Document 259: total_number: 1, positions: [64]
Document 300: total_number: 2, positions: [50, 55]
Document 304: total_number: 2, positions: [150, 244]
Document 306: total_number: 1, positions: [163]
Document 314: total_number: 8, positions: [54, 58, 73, 77, 98, 154, 158, 172]
Document 319: total_number: 1, positions: [198]
Document 353: total_number: 1, positions: [33]
Document 354: total_number: 1, positions: [338]
```

### زیر بخش 3:

### پرسمان 1:

```
Query: باشگاه های فوتبال آسیا
Terms: ['باشگاه', 'آسیا', 'فوتسال']
Document: 0 in Title of: باشگاه های فوتبال آسیا Relevance Score: 3
Document: 797 in Title of: Relevance Score: 3 پتانسیل تهران برای میزبانی جام باشگاه‌ها بیشتر است/حافظیه در حد لیگ یک و دو هم نیست
Document: 1306 in Title of: Relevance Score: 3 تفریح جزئیات دریافت مجوز حرفه‌ای برای تیم‌های فوتبال و تشکیل کمیته ویژه در هفته‌های پایانی لیگ
Document: 1535 in Title of: Relevance Score: 3 چرا تیم ملی فوتبال به تایلند اعزام نمی‌شود؟
Document: 1671 in Title of: Relevance Score: 3 از فوتبال به فوتبال رسید/باشگاه‌های ایران به هوش باشند AFC مجوز حرفه‌ای
```

برای مثال در عنوان داک شماره 0، می‌توان عبارت سرچ شده را مشاهده کرد.

### پرسمان 2:

```
Query: باشگاه های فوتبال ! آسیا
Terms: ['باشگاه', 'فوتسال']
Forbidden Terms: ['آسیا']
phase words: {}
Document: 81 in Title of: Relevance Score: 2 ماجدی: فوتبال کشور به تغییرات نیاز دارد
Document: 426 in Title of: Relevance Score: 2 اتفاق عجیب در لیگ نوجوانان/ قهرمانی سایپا پس گرفته شد!نامه
Document: 453 in Title of: Relevance Score: 2 دبیر هیات فوتبال اصفهان: ابلاغی برای حضور تماشاگران نداشتیم/ امیدوارم روزی مصوات رعایت شود
Document: 467 in Title of: Relevance Score: 2 سرمربی شهرخودرو جرعه شد/محرومیت برای بازیکن فوتبال
Document: 493 in Title of: Relevance Score: 2 اوج حساسیت در هفته آخر لیگ برتر فوتبال/ انصراف شهروند از بازی مقابل مس
Document: 1339 in Title of: Relevance Score: 2 فولاد در کشتی تیمداری می‌کند/گرشاسمی: در جذب بازیکن رکورد زده ایم
Document: 1594 in Title of: Relevance Score: 2 مستند/مجاز زندی بتن را کمیته داوران تایید کرد AFC جابری: همه باشگاه‌ها ملزم به رعایت ابلاغیه
Document: 1780 in Title of: Relevance Score: 2 خوراکی: داور دو گل صحیح ما را مردود اعلام کرد/کراپ مظلوم واقع شد
Document: 2033 in Title of: Relevance Score: 2 شکایت باشگاه گیتی پسند از سرمربی مس سونگون ورزشان
Document: 2320 in Title of: Relevance Score: 2 بازیکنان فوتبال مس برای راستی آزمایی تست کرونا به تهران می‌آیند/ مصاف گیتی پسند و مس 29 بهمن
Document: 2391 in Title of: Relevance Score: 2 بازیکن سایر نگرانان می‌توانند در لیگ امید بازی کنند/ پیکان در بودجه ضعیف دارد
Document: 2435 in Title of: Relevance Score: 2 حضور مدیرعامل جدید در تمرین پیکان و وعده به شاگردان حمینی+عکس
Document: 2453 in Title of: Relevance Score: 2 تصمیمات سازمان لیگ فوتبال راهگشا یا بلا جان؟
Document: 2478 in Title of: Relevance Score: 2 بیانیه مس سونگون در خصوص تصمیم سازمان لیگ فوتبال
Document: 2583 in Title of: Relevance Score: 2 اقدام ارزشمند گیتی پسند در پی مثبت شدن کرونای بازیکنان مس سونگون
```

برای جلوگیری از نمایش داک های شامل کلمه آسیا، مقدار score آن را 1- می‌کنیم. برای مثال داک شماره 0 را در زیر مشاهده می‌کنید:

```
Document: 6763 in Title of: Relevance Score: 1 جابری: نامه ای درباره انصراف شهروند از حضور در لیگ برتر فوتبال ارسال ندهد است
Document: 9712 in Title of: Relevance Score: 1 ناکید رئیسجمهور در تدوین بودجه ۱۴۰۱ بر عدالت بود/ عزم جدی دولت برای حل مشکلات
Document: 0 in Title of: Relevance Score: -1 اعلام زمان قرعه کشی جام باشگاه های فوتبال آسیا
Document: 8 in Title of: Relevance Score: -1 مدیر بی‌ستاورد، رئیس شد/ روزگار تلختر از تلخ برای والیبال ارومیه؟
Document: 29 in Title of: Relevance Score: -1 مدیرعامل آئومینیموم: مرفاوی می‌داند هیچ باشگاهی با یک نفر صحبت نمی‌کند/استقلال‌ها لیگ را گرفته‌اند
Document: 50 in Title of: Relevance Score: -1 یک شاکای دیگر به والیبال ارومیه اضافه شد/مجرد: قرارداد را گم، پشتم را خالی و بلاکم کردند
Document: 54 in Title of: Relevance Score: -1 اکبری: سیاست پیکان کسب قهرمانی نبود/شهادت بزد برای قهرمانی بسته شده بود
```

### پرسمان 3:

```
Query: 'سهمیه المپیک'
Terms: {}
Forbidden Terms: []
phase words: [['سهمیه', 'المپیک']]
Document: 3907 in Title of: Relevance Score: 1 تیراندازان در کدام رویدادها المپیکی می‌شوند؟
```

## پرسمان 4:

در این پرسمان عبارت “طلای ‘لیگ برتر’ چغر” را سرچ کرده ام:

```
Query: 'لیگ برتر'
Terms: {}
Forbidden Terms: []
phase words: [['لیگ', 'برتر']]
Document: 1243 in Title of: هفته بیستم لیگ برتر/ استقلال 20 می شود؟ / پرسپولیس مقابل دومین تیم کرمانی/دوئل مریدان استقلال Relevance Score: 2
Document: 3518 in Title of: هفته بیستم لیگ برتر فوتسال/ صعود سن ایچ به رده دوم با شکست سیاهان Relevance Score: 2
Document: 5294 in Title of: هفته چهاردهم لیگ برتر/کورس سרחابی‌ها و سیاهان برای قهرمانی نیم فصل/استقلال مقابل حریفی چغر، دیدار پرسپولیس با قهرنشین Relevance Score: 2
Document: 2 in Title of: محل برگزاری نقشه‌های خبری سרחابی‌ها؛ مجیدی در سازمان لیگ، گل‌محمدی در تمرین پرسپولیس Relevance Score: 1
Document: 261 in Title of: هفته پایانی، لیگ برتر فوتسال | گیت، یسند با شکست جیسم قهرمانی، خود را در اصفهان جشن گرفت Relevance Score: 1
```

```
Query: طلای 'لیگ برتر'
Terms: ['طلا']
Forbidden Terms: []
phase words: [['لیگ', 'برتر']]
Document: 974 in Title of: روایتی جالب از تیم قهرمان لیگ برتر واترپلو/ به نام شهیدی که تیم ملی و مدال را رها کرد و رفت Relevance Score: 2
Document: 1243 in Title of: هفته بیستم لیگ برتر/ استقلال 20 می شود؟ / پرسپولیس مقابل دومین تیم کرمانی/دوئل مریدان استقلال Relevance Score: 2
Document: 3518 in Title of: هفته بیستم لیگ برتر فوتسال/ صعود سن ایچ به رده دوم با شکست سیاهان Relevance Score: 2
Document: 5294 in Title of: هفته چهاردهم لیگ برتر/کورس سרחابی‌ها و سیاهان برای قهرمانی نیم فصل/استقلال مقابل حریفی چغر، دیدار پرسپولیس با قهرنشین Relevance Score: 2
Document: 2 in Title of: محل برگزاری نقشه‌های خبری سרחابی‌ها؛ مجیدی در سازمان لیگ، گل‌محمدی در تمرین پرسپولیس Relevance Score: 1
```

```
Query: طلای 'لیگ برتر' ! چغر
Terms: ['طلا']
Forbidden Terms: ['چغر']
phase words: [['لیگ', 'برتر']]
Document: 974 in Title of: روایتی جالب از تیم قهرمان لیگ برتر واترپلو/ به نام شهیدی که تیم ملی و مدال را رها کرد و رفت Relevance Score: 2
Document: 1243 in Title of: هفته بیستم لیگ برتر/ استقلال 20 می شود؟ / پرسپولیس مقابل دومین تیم کرمانی/دوئل مریدان استقلال Relevance Score: 2
Document: 3518 in Title of: هفته بیستم لیگ برتر فوتسال/ صعود سن ایچ به رده دوم با شکست سیاهان Relevance Score: 2
Document: 2 in Title of: محل برگزاری نقشه‌های خبری سרחابی‌ها؛ مجیدی در سازمان لیگ، گل‌محمدی در تمرین پرسپولیس Relevance Score: 1
Document: 261 in Title of: هفته پایانی، لیگ برتر فوتسال | گیت، یسند با شکست جیسم قهرمانی، خود را در اصفهان جشن گرفت Relevance Score: 1
```

همان‌طور که مشاهده می‌کنید با اضافه کردن کلمه طلا به عبارت لیگ برتر، داک شماره 974 که قبلا امتیاز 1 داشت، نیز امتیاز 2 را بدست آورد. در ادامه با ممنوع کردن کلمه چغر، داک شماره 5294 از دور رقابت حذف گردید.

## پرسمان 5:

```
Query: 'مایکل ! جردن'
Terms: {}
Forbidden Terms: ['جردن']
phase words: [['مایکل', 'جردن']]
No result for this query
```

در ادامه تعداد تکرار کلمات را نیز به صورت یک امتیاز مجزا در نظر می‌گیریم و اولویت را ابتدا دارا بودن تمام کلمات و سپس تکرار کلمات در نظر می‌گیریم:

برای مثال، برای کوئری “ ‘لیگ برتر’ چغر ” خواهیم داشت:

```
Query: 'لیگ برتر' ! چمر
Terms: []
Forbidden Terms: ['چمر']
phase words: [['لیگ', 'برتر']]
Document: 1243 in Title of: هفته بیستم لیگ برتر| استقلال 20 می شود؟ / پرسپولیس مقابل دومین تیم کرمانی/دوئل مرینیان استقلال Relevance Score: (2, 46)
Document: 3518 in Title of: هفته بیستم لیگ برتر فوتبال| صعود سن ایچ به رده دوم با شکست سیاهان Relevance Score: (2, 4)
Document: 5918 in Title of: اجرای تصمیم گیری در والیبال/سیاهان از بین دو گزینه، سومی را انتخاب کرد Relevance Score: (1, 37)
Document: 1634 in Title of: هفته نوزدهم لیگ برتر| استقلال بی دفاع مقابل فجر؛ دوئل قلعه نویی با پرسپولیس و دربی اصفهان Relevance Score: (1, 27)
Document: 4624 in Title of: نکته از لیگ برتر در نیم فصل اول| پرسپولیس و 5 تیم رکورددار حاشیه / زمین مفت نوبکتیا و لطف قلعه نویی به استقلال 16 Relevance Score: (1, 24)
Document: 6673 in Title of: ایساج مدیر نظارت و یازوسی سازمان لیگ؛ آقای اصولی چه آورده ای برای فوتبال خراسان و پدیده داشتی؟ Relevance Score: (1, 24)
Document: 820 in Title of: یکم لیگ برتر|تصوریان به دنبال عاقبت بخیری با پرسپولیس/تقابل استقلال و قلعه نویی بدون بازیکن گابیتی Relevance Score: (1, 23)
Document: 5632 in Title of: هفته سیزدهم لیگ برتر فوتبال| تداوم شکست ناپذیری استقلال در اراک؟/ نبرد قهرمزبوشان تهران و تبریز Relevance Score: (1, 23)
Document: 5992 in Title of: هفته دوازدهم لیگ برتر|سیاهان - پرسپولیس؛ نبرد آسیایی در زمین بیطرف/استقلال مقابل مانع فولادی و دربی کرمان Relevance Score: (1, 23)
```