



۱۴۰۰/۰۹/۰۱

# بازیابی هوشمند اطلاعات تمرین اول





در ابتدا محیط لازم برای انجام تمرین را به این ترتیب فراهم کردیم:

- دانلود و نصب نسخه ۲۰.۴ از سیستم عامل اوبونتو
- نصب جاوا
- نصب Maven
- نصب Galago

## پیش‌نیاز: ایجاد شاخص

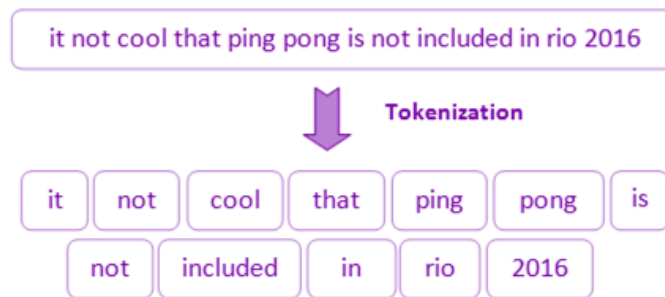
در هر زبانی، کلمات با توجه به نقشی که در جملات ایفا می کنند، به شکل‌های ظاهری متفاوتی خواهند بود. اما تمامی آن‌ها از یک ریشه ساخته می شوند. لذا در بسیاری از روش‌ها، ابتدا می‌بایست ریشه کلمات را پیدا کنیم. یکی از روش‌های متداول برای ریشه‌یابی کلمات، روش ریشه‌یابی Stemming است. الگوریتم‌های مختلفی جهت انجام عمل ریشه‌یابی وجود دارد که الگوریتم Porter از الگوریتم‌های معروف در زبان انگلیسی می باشد. این الگوریتم طبق یک سری قاعده‌ی منظم (مثلاً حذف حرف s در آخر کلمات جمع) می‌تواند ریشه‌ی کلمات را با دقت خوبی به دست آورد.

پس از نصب و راه‌اندازی گالاگو، به کمک دستورات گالاگو از روش Porter Stemmer برای ریشه‌یابی کلمات استفاده کردیم.

```
"stemmer" : ["porter"],
```

عمل Tokenization متن را مانند شکل زیر به توکن‌های تشکیل دهنده خود تبدیل می‌کند. این عمل را با دستورات زیر انجام می‌دهیم:

```
"tokenizer" : {  
  "fields" : ["text", "head"],  
  "formats" : {  
    "text" : "string",  
    "head" : "string"  
  }  
}
```



نهایتاً فرمت فایل که در ابتدا به صورت بدون فرمت بود را با کمک نرم افزار zip7 از سیستم عامل ویندوز اکسترکت کرده و به فایل txt. رسیدیم، سپس فایل حاصل شده را به فرمت trectext. که قالبی مشابه XML دارد بردیم. این تایپ از فایل ها برای اسناد و متونی مناسب است که لازم داریم قسمت های مختلف آن را بتوان جدا دید.

```
<DOC>
<DOCNO> AP890325-0001 </DOCNO>
<FILEID>AP-NR-03-25-89 0106EST</FILEID>
<FIRST>r a AM-People-Bridges 1stLd-Writethru a0733 03-25 0278</FIRST>
<SECOND>AM-People-Bridges, 1st Ld-Writethru, a0733,0282</SECOND>
<HEAD>Bridges Pleads Innocent To Attempted Murder Charge</HEAD>
<HEAD>Eds: SUBS lead to reflect that Bridges' role in the shooting has
not been established.</HEAD>
<DATELINE>LOS ANGELES (AP) </DATELINE>
<TEXT>
Actor Todd Bridges, a star in the NBC comedy
'Diff'rent Strokes,' pleaded innocent Friday to attempted murder
charges stemming from an incident in which he allegedly shot a
roommate.
Bridges shot Kenneth Clay as many as eight times in a Feb. 2
argument that stemmed from Clay having borrowed Bridges' BMW
automobile, according to testimony in the case.
Another man, Harvey Duckett, 30, who was with Bridges during the
episode, also faces the attempted murder charge. But Duckett, who
has pleaded no contest, is testifying in the case in exchange for
leniency from prosecutors.
Bridges, 23, was bound over to Superior Court last week by
Municipal Judge David Horwitz, who also refused to reduce Bridges'
$2 million bail.
On Friday, Superior Court Judge David Horowitz scheduled a
pretrial hearing, trial setting and bail motions for April 14.
In addition to his role as Gary Coleman's protective older
brother on 'Diff'rent Strokes,' Bridges has appeared in ABC's
'Fish' and has been on various installment of 'Circus of the
Stars.'
During a preliminary hearing earlier this month, Bridges'
attorney, Johnnie Cochran argued for a reduction in charges,
contending Bridges was too drugged to premeditate attempted
first-degree murder. Witnesses testified that the actor ingested
rock cocaine at least four times on the day of the shooting.
Clay described Bridges as 'based out' from freebasing or
smoking cocaine. 'It was the worst I seen him. He looked like his
eyes were about to jump out of his head,' Clay testified.
</TEXT>
</DOC>
```

دستورات مربوط به تنظیماتی که در بالا بدان ها اشاره شد را در فایل indexSettings.json قرار داده و

دستور زیر را در ترمینال اجرا می کنیم:

```
Galago/galago-3.16/core/target/appassembler/bin/galago build
/home/arya/Desktop/CA1-Resources/indexSettings.json
```

پس از گذشت حدوداً یک ساعت از اجرای دستور Build عمل شاخص گذاری با موفقیت به اتمام رسید.

Done Indexing.

- 0.92 Hours
- 55.18 Minutes



- 3310.79 Seconds  
Documents Indexed: 163912.

## سوال ۱: تابع بازیابی BM25

**الف)** در این بخش با روش BM25 بازیابی را ابتدا با مقادیر پیش فرض  $b$  و  $k$  روی مجموعه پرس و جو ۱۰۱ الی ۱۵۰ با مقدار ۱۰۰ برای تعداد Requested، انجام داده و مقادیر Recall، MAP، nDCG و  $P@5$  را بدست می آوریم. در ادامه مقادیر  $b$  و  $k$  را به صورت آزمون و خطا ابتدا با گام های بلند تغییر داده و هر نوبت مقادیر را یادداشت می کنیم و در صورتی که شاهد بهبود نسبت به حالت پیش فرض بودیم، مقادیر نزدیک را با گام های کوچک تری امتحان می کنیم تا به مقدار بهینه برسیم.

در جدول زیر نتایج بدست آمده قابل مشاهده است:

۱۰۱-۱۵۰						k	b	
P@5	nDCG	MAP	Recall					
۰.۳۶۸	۰.۳۳۱	۰.۱۶۸	۰.۲۲۹	۴۸۰.۵	۱۰۹۹	۱.۲	۰.۷۵	BM25 روش ۱
۰.۳۶۸	۰.۳۱۸	۰.۱۵۹	۰.۲۲۱	۴۸۰.۵	۱۰۶۱	۱.۲	۰.۳۰	
۰.۳۷۶	۰.۳۳۲	۰.۱۶۵	۰.۲۲۹	۴۸۰.۵	۱۱۰۰	۱.۲	۰.۷۰	
۰.۳۸۰	۰.۳۳۵	۰.۱۶۹	۰.۲۳۴	۴۸۰.۵	۱۱۲۴	۱.۲	۰.۵۵	
۰.۳۸۰	۰.۳۳۵	۰.۱۶۹	۰.۲۳۴	۴۸۰.۵	۱۱۲۴	۱.۲	۰.۵۵	
۰.۴۱۲	۰.۳۴۱	۰.۱۷۴	۰.۲۴۰	۴۸۰.۵	۱۱۵۱	۱.۵	۰.۵۰	
۰.۴۰۸	۰.۳۴۵	۰.۱۷۶	۰.۲۴۱	۴۸۰.۵	۱۱۵۶	۱.۷	۰.۵۰	
۰.۳۹۲	۰.۳۴۴	۰.۱۷۳	۰.۲۳۸	۴۸۰.۵	۱۱۴۵	۱.۷	۰.۶۰	
۰.۴۰۸	۰.۳۴۵	۰.۱۷۷	۰.۲۴۰	۴۸۰.۵	۱۱۵۳	۱.۸	۰.۵۰	
۰.۴۱۶	۰.۳۴۷	۰.۱۷۸	۰.۲۴۰	۴۸۰.۵	۱۱۵۴	۲.۰	۰.۵۰	
۰.۴۰۰	۰.۳۴۶	۰.۱۷۷	۰.۲۳۶	۴۸۰.۵	۱۱۳۵	۲.۵	۰.۶۰	
۰.۴۲۰	۰.۳۴۸	۰.۱۸۲	۰.۲۴۳	۴۸۰.۵	۱۱۶۸	۲.۶	۰.۴۰	
۰.۳۹۲	۰.۳۴۲	۰.۱۷۳	۰.۲۳۰	۴۸۰.۵	۱۱۰۷	۲.۵	۰.۷۰	
۰.۳۹۲	۰.۳۴۴	۰.۱۷۹	۰.۲۴۳	۴۸۰.۵	۱۱۶۹	۱.۹	۰.۳۰	
۰.۳۴۰	۰.۲۲۹	۰.۱۴۳	۰.۱۹۷	۴۸۰.۵	۹۴۵	۲.۰	۱.۰۰	

با جست‌وجو در منابع علمی محدوده مناسب برای پارامتر  $b$ ،  $0.3$  الی  $1$  و محدوده مناسب برای پارامتر  $k$ ،  $0.5$  الی  $2.5$  به نظر آمد و بر این اساس به تست کردن مقادیر در این نواحی و نیز حاشیه‌ای از بالا و پایین این نواحی پرداختیم و مقایسه بهینه برای پرس‌وجوهای  $101$  الی  $150$  برای پارامترهای  $b$  و  $k$  به ترتیب  $0.4$  و  $2.6$  دیده شد.

نتایج حاصل از این مقادیر به نسبت نتایج بدست آمده از مقادیر پیش‌فرض (ردیف اول جدول) مقدار بیشتر داشت که نشان‌دهنده این است که بهینه‌سازی این پارامترها در کسب نتایج بهتر موفق بوده است.

**ب)** در این بخش مقادیری که برای پارامترهای  $b$  و  $k$  روی پرس‌وجوهای  $101$  الی  $150$  بهینه بدست آوردیم را روی پرس‌وجوهای  $51$  الی  $100$  اعمال کرده و با مقادیر پیش‌فرض مورد مقایسه قرار می‌دهیم.

۵۱-۱۰۰						k	b	
P@5	nDCG	MAP	Recall					
۰.۳۸۲	۰.۳۲۲	۰.۱۷۰	۰.۲۲۰	۱۰۶۲۰	۲۳۳۹	۱.۲	۰.۷۵	۱-روش BM25
۰.۴۱۲	۰.۳۳۱	۰.۱۷۹	۰.۲۲۸	۱۰۶۲۰	۲۴۱۸	۲.۶	۰.۴۰	

همان‌طور که ملاحظه می‌شود با استفاده از مقادیر بهینه برای پارامترهایی که در قسمت قبل برای پرس‌وجوهای متفاوتی بدست آوردیم، روی این مجموعه از پرس‌وجو هم نتیجه بهتری از مقادیر پیش‌فرض رسیده ایم. این مسئله نشان دهنده اثر مجموعه اسناد روی بهبود مقدارهای پارامترها می‌باشد.

در ادامه از مقادیری که به عنوان مقادیر بهینه برای پارامترها محاسبه کردیم و نیز مقدار پیش‌فرض و دومین بهترین مقدار بهینه‌ای که در بخش قبلی استخراج کرده‌ایم استفاده کرده و روش‌های پیشنهاد شده در تمرین را پس از پیاده‌سازی برای پرس‌وجوهای  $51$  الی  $100$  مورد بازیابی و ارزیابی قرار می‌دهیم. نتایج حاصل شده در جدول زیر ارائه شده.

۵۱-۱۰۰						k	b	
P@5	nDCG	MAP	Recall					
۰.۳۸۲	۰.۳۲۲	۰.۱۷۰	۰.۲۲۰	۱۰۶۲۰	۲۳۳۹	۱.۲	۰.۷۵	۱-روش BM25
۰.۴۱۲	۰.۳۳۱	۰.۱۷۹	۰.۲۲۸	۱۰۶۲۰	۲۴۱۸	۲.۶	۰.۴۰	
۰.۰۲۲	۰.۰۲۱	۰.۰۰۸	۰.۰۱۷	۱۰۶۲۰	۱۷۸	۱.۲	۰.۷۵	۲-روش اول
						۲.۶	۰.۴۰	
						۱.۹	۰.۳۰	
۰.۴۰۲	۰.۲۸۷	۰.۱۴۷	۰.۲۰۱	۱۰۶۲۰	۲۱۳۰	۱.۲	۰.۷۵	۳-روش دوم



۰.۳۹۲	۰.۲۷۹	۰.۱۴۱	۰.۱۹۴	۱۰۶۲۰	۲۰۵۹	۲.۶	۰.۴۰	
۰.۴۰۴	۰.۲۸۴	۰.۱۴۵	۰.۱۹۹	۱۰۶۲۰	۲۱۱۴	۱.۹	۰.۳۰	
						۱.۲	۰.۷۵	
۰.۲۲۲	۰.۱۸۱	۰.۰۷۳	۰.۱۲۶	۱۰۶۲۰	۱۳۳۶	۲.۶	۰.۴۰	۴-روش سوم
						۱.۹	۰.۳۰	
۰.۳۹۶	۰.۳۲۴	۰.۱۷۱	۰.۲۲۰	۱۰۶۲۰	۲۳۴۱	۱.۲	۰.۷۵	
۰.۴۱۶	۰.۳۲۹	۰.۱۷۸	۰.۲۲۶	۱۰۶۲۰	۲۴۰۵	۲.۶	۰.۴۰	۵-روش چهارم
۰.۴۱۲	۰.۳۳۰	۰.۱۷۹	۰.۲۲۹	۱۰۶۲۰	۲۴۳۳	۱.۹	۰.۳۰	
۰.۳۸۲	۰.۳۲۲	۰.۱۷۰	۰.۲۲۰	۱۰۶۲۰	۲۳۳۹	۱.۲	۰.۷۵	
۰.۴۱۲	۰.۳۳۱	۰.۱۷۹	۰.۲۲۸	۱۰۶۲۰	۲۴۱۹	۲.۶	۰.۴۰	۶-روش پنجم با
۰.۴۰۶	۰.۳۳۲	۰.۱۸۰	۰.۲۳۰	۱۰۶۲۰	۲۴۴۳	۱.۹	۰.۳۰	مقدار تتای ۱

فایل‌های پیاده‌سازی‌های مربوط به روش‌های پیشنهادی اول تا پنجم در فایل‌های av2 الی av6 در کنار گزارش قرار دارند که به دلیل این که در صورت تمرین اشاره شده که به توضیح آن‌ها نپردازیم از نوشتن در خصوصشان صرف نظر کردیم.

در تفسیر نتایج بدست آمده که در جدول بالا ارائه شده‌اند باید اشاره کرد که روش پیشنهادی اول عملکرد بسیار ضعیفی از خود نشان داده و چون در آن از پارامترهای  $b$  و  $k$  استفاده نشده است، لذا بهینه‌سازی آن‌ها اثر در بهبود یافتن عملکرد بازیابی ندارد. همچنین نتایج بدست آمده از این روش عملکرد بسیار ضعیفی نسبت به سایر روش‌های پیشنهادی و نیز روش اصلی دارد.

روش پیشنهادی دوم پس از لحاظ کردن پارامترهای بهینه به مقدار کمتری برای معیار MAP نسبت به استفاده از مقادیر پیش‌فرض دست پیدا کرد اما با این حال از عملکرد قابل قبولی از خود نشان داد.

روش پیشنهادی سوم علی‌رغم سادگی بسیار عملکردی به مراتب مناسب‌تر از روش اول بر روی داده‌ها و پرس‌وجوهای ارائه شده از خود به نمایش گذاشت هرچند که در کل نتایج جالبی نسبت به روش‌های بعدی نداشت. در این روش هم به دلیل این که در آن از پارامترهای  $b$  و  $k$  استفاده نشده است، عملیات بهینه‌سازی‌ای که در بخش‌های قبل روی پارامترها داشتیم برای بهبود دادن عملکرد بی‌اثر بود.



در روش پیشنهادی چهارم به نتایج بهتری نسبت به روش‌های پیشنهادی قبلی و نیز روش اصلی دست پیدا کردیم. در این روش هم شاهد بهبود قابل ملاحظه معیارهای ارزیابی برای استفاده از پارامترهای بهینه یافته شده در قسمت‌های قبل نسبت به مقادیر پیش فرض این پارامترها هستیم.

در روش پیشنهادی پنجم از ما خواسته شده که فورمول ارائه شده را پیاده‌سازی کرده و برای مقادیر مختلف از تتا مورد بررسی قرار دهیم و در نهایت نتایج بدست آمده با تتای برتر را گزارش نمائیم. با بررسی‌های صورت گرفته بهترین مقدار برای تتا در این روش ۱ اندازه‌گیری شده که مقادیر معیارهای ارزیابی مرتبط به آن در جدول قرار دارند. در این روش با بهینه‌سازی مقدار پارامترها موفق شدیم که به نتایج بهتری برای همه معیارهای ارزیابی نسبت به حالت اولیه پارامترها دست پیدا کنیم.

در انتها این بخش جهت سهولت دسترسی بخش‌هایی از قسمت‌های مهم‌تر از کد ضمیمه شده است.

```
private double score(double count, double length) {  
    return idf;  
}
```

```
private double score(double count, double length) {  
    double numerator = count * (k + 1);  
    double denominator = count + k;  
    return numerator / denominator;  
}
```

```
private double score(double count, double length) {  
    double i = 0;  
    if(count != 0){  
        i = 1;  
    }  
    return i;  
}
```

```
private double score(double count, double length) {  
    double numerator = (k + 1) * (count / (1 - b + (b * length /  
avgDocLength)) + 0.5);  
    double denominator = k + (count / (1 - b + (b * length /  
avgDocLength)) + 0.5);  
    return idf * numerator / denominator;  
}
```

```
private double score(double count, double length) {  
    double numerator = count * (k + 1);
```



```
double denominator = count + (k * (1 - b + (b * length /
avgDocLength)));
return idf * (numerator / denominator + 1);
}
```

```
// count -> score iterators
// Scorers can be named directly as nodes
{av2.class.getName(), "av2"},
{av3.class.getName(), "av3"},
{av4.class.getName(), "av4"},
{av5.class.getName(), "av5"},
{av6.class.getName(), "av6"},
```

## سوال ۲: تابع بازیابی Pivoted Length Normalization

در جدول زیر مقایسه روش‌های خواسته شده در صورت تمرین برای معیارهای کارایی روی پرس‌جوهای ۵۱-۱۰۰ قابل مشاهده است. با توجه به جدول زیر روش‌های BM25 و BM25+ از نظر معیار کارایی MAP نتایج بهتری برای مقادیر پیش‌فرض پارامترهای b و k دارند.

۵۱-۱۰۰						
P@5	nDCG	MAP	Recall			
۰.۲۰۳	۰.۲۸۵	۰.۱۲۱	۰.۲۰۶	۱۰۶۲۰	۲۱۸۸	مدل اصلی
۰.۳۸۲	۰.۳۲۲	۰.۱۷۰	۰.۲۲۰	۱۰۶۲۰	۲۳۳۹	BM25
۰.۳۸۲	۰.۳۲۲	۰.۱۷۰	۰.۲۲۰	۱۰۶۲۰	۲۳۳۹	BM25+
۰.۲۰۱	۰.۲۶۱	۰.۱۱۹	۰.۱۹۴	۱۰۶۲۰	۲۰۶۸	مدل بدون مؤلفه لگاریتمی تو در تو

در انتها جهت سهولت دسترسی بخش‌هایی از قسمت‌های مهم‌تر از کد ضمیمه شده است.

```
private double score(double count, double length) {
double numerator = log(1 + log(1 + count));
double denominator = 1 - b + (b * length / avgDocLength);
log = log((documentCount + 1) / (df));
```





```
return count * (numerator / denominator) * (logg);  
}
```

```
private double score(double count, double length) {  
    double numerator=1 + log(1 + count);  
    double denominator = 1 - b + (b * length / avgDocLength);  
    logg = log(( documentCount + 1) / (df));  
    return count*(numerator/denominator)*(logg);  
}
```

```
{avb1.class.getName(), "avb1"},  
{avb2.class.getName(), "avb2"},
```