

گزارش مرحله دوم فاز اول پروژه درس بازیابی اطلاعات – ساخت شاخص مکانی

نگار موقتیان، ۹۸۳۱۰۶۲

نحوه ساخت شاخص مکانی را با ذکر نمونه خروجی نمایش دهید.

برای ساخت شاخص مکانی در ابتدا چند کلاس تعریف شده‌اند تا بتوانند اطلاعات مربوط به قسمت‌های مختلف شاخص را نگهداری کنند.

ساختمان داده کلی استفاده شده dictionary است که کلمات (term) مختلف را به یک شیء از کلاس Index نگاشت می‌کند (مزیت استفاده از dictionary در سرعت بالای دسترسی به value مربوط به هر یک از term هاست). کلاس Index خود تعداد اسنادی که کلمه فوق در آن‌ها تکرار شده (doc. frequency) و همچنین یک لیست از اشیاء کلاس Postings را نگهداری می‌کند (postings list). کلاس Postings نیز یک شماره سند دارد که نشان می‌دهد کلمه فوق در آن سند چند بار (term frequency) و در چه مکان‌هایی تکرار شده.

اولین قدم برای ساخت شاخص مکانی مرتب‌سازی توکن‌های بدست آمده از مرحله قبل بر اساس حروف الفباست تا بتوانیم کلمات یکسان را یکی کرده و شاخص مکانی هر یک را بسازیم. پس از آن به ترتیب بر روی تمام توکن‌ها پیمایش کرده و مانند زیر عمل می‌کنیم:

۱. یک Index جدید بساز

۲. تا وقتی که ترم توکن‌ها یکسان است:

۳. یک Postings جدید بساز

۴. تا وقتی که شماره اسناد یکسان است:

۵. توکن را به Postings اضافه کن

۶. Postings را به Index اضافه کن

۷. ترم را به همراه Index به دیکشنری شاخص مکانی اضافه کن

با این کار همان ساختاری میان کلاس‌ها که توضیح داده شد رعایت خواهد شد. لیست مکان کلمات یکسان در اسناد یکسان درون یک Postings مشخص قرار خواهد گرفت و لیست اسناد با کلمات یکسان نیز درون یک Index مشخص قرار خواهد گرفت. در نهایت نیز به تدریج با افزودن کلمات متفاوت دیکشنری ساخته خواهد شد.

در نهایت خروجی برنامه به ازای توکن‌های مرتب‌شده زیر:

```
sorted tokens: [0 - 6 : AFC, 3 - 146 : یآ, 3 - 44 : مآرآ, 3 - 31 : ووزرآ, 2 - 50 : دآزآ, 0 - 5 : ایسآ, 0 - 21 : ایسآ, 0 - 35 : ایسآ, 0 - 57 : ایسآ, 3 - 49 : وروآ#دروآ, 0 - 65 : نهآ, 4 - 40 : یگدامآ, 3 - 18 : هدامآ, 3 - 145 : هدامآ, 3 - 67 : هدامآ, 3 - 46 : هدامآ, 3 - 28 : یسآ, 3 - 131 : یسآ, 3 - 27 : 1, راهظا, 1 - 17 : عاظا, 1 - 20 : دنقسا, 2 - 24 : لالقسا, 2 - 30 : لالقسا, 2 - 8 : ساسا, 0 - 24 : دنمشرآ, 3 - 102 : ودرا, 3 - 144 : قالخا, 3 - 36 : 3 - 69 : کیپملا, 3 - 42 : کیپملا, 3 - 23 : کیپملا, 3 - 10 : مالعآ, 3 - 50 : مالعآ, 0 - 23 : افعآ, 3 - 98 : افعآ, 3 - 79 : مازعآ, 3 - 128 : راهظا, 3 - 27 : کیپملا, 3 - 148 : کیپملا, 3 - 142 : کیپملا, 3 - 109 : کیپملا, 3 - 90 : کیپملا]
```

مانند زیر می‌باشد:

```
positional index: {'AFC': #1 -> [0 : #1 -> [6]], 'یآ': #1 -> [3 : #1 -> [146]], 'مآرآ': #1 -> [3 : #1 -> [44]], 'ووزرآ': #1 -> [3 : #1 -> [31]], 'دآزآ': #1 -> [2 : #1 -> [50]], 'ایسآ': #1 -> [0 : #4 -> [5, 21, 35, 57]], 'یسآ': #1 -> [3 : #2 -> [49, 131]], 'هدامآ': #1 -> [3 : #4 -> [28, 46, 67, 145]], 'نهآ': #1 -> [4 : #1 -> [40]], 'یگدامآ': #1 -> [3 : #1 -> [18]], 'دروآ': #1 -> [0 : #1 -> [65]], 'عاظا': #1 -> [3 : #1 -> [36]], 'دنقسا': #1 -> [2 : #2 -> [8, 30]], 'لالقسا': #1 -> [2 : #1 -> [24]], 'ساسا': #1 -> [0 : #1 -> [24]], 'دنمشرآ': #1 -> [3 : #1 -> [102]], 'ودرا': #1 -> [3 : #1 -> [144]], 'قالخا': #1 -> [1 : #1 -> [20]], 'راهظا': #1 -> [1 : #2 -> [17, 27]], 'مآزعا': #1 -> [3 : #1 -> [128]], 'افعا': #1 -> [3 : #2 -> [79, 98]], 'مالعآ': #2 -> [0 : #1 -> [23]], 3 : #1 -> [50]], 'کیپملا': #1 -> [3 : #8 -> [10, 23, 42, 69, 90, 109, 142, 148]]}
```

* در مثال بالا تنها بخشی از توکن‌ها از اسناد ۱ تا ۵ آورده شده‌اند.