محمد مهدی نظری - ۹۹۳۱۰۶۱

در ساخت شاخص مکانی از لینک زیر کمک گرفته شده:

https://www.geeksforgeeks.org/python-positional-index/

برای نگه داری توکن ها از یک دیکشنری استفاده میکنیم که key آن همان توکن و value آن لیستی از تعداد تکرارهای آن در کل داکیومنت ها و دیکشنری از موقعیت تکرار آن در هر داکیومنت است : دیکشنری دوم کلید برابر DocID و مقدار برابر لیستی از پوزیشن هایی است که آن توکن در داکیومنت آمده است(برای هر داکیومنت یک دیکشنری موجود است).

مراحل ساخت Positional Index در کد بصورت خط به خط با کامنت توضیح داده شده است. تابع فعلی فقط برای یک فایل تکست نوشته شده پس فقط یک داکیومنت داریم که ID آن برابر ۰ در نظر گرفته شده اما در گام بعدی پروژه، تابع برای کار با فایل json بروزرسانی خواهد شد و داکیومنت های مختلف بررسی میشوند.

برای استفاده ازتابع ابتدا باید متن را پیش پردازش کنیم یعنی نرمال سازی و حذف کلمات پرتکرار و ریشه یابی انجتم شده و توکن ها استخراج میشوند و متن نهایی به عنوان آرگومان به تابع تحویل داده میشود. خروجی برنامه را با متن ورودی

"بارها و بارها از نخبگان و دلسوزان نظام این جمله بیان شده که قبل از وقوع انقلاب بلکه از هزاران سال پیش مردم ایران به پوشش و حجاب ارج مینهادند و در حال حاضر این مسئله همچون رفتار رضا شاه میشه .خوبه که داریم می رویم به سمت حجاب اختیاری و نظام مورد علاقه ما مثل رضا شاه نمی شه".

ملاحظه میکنیم:

"C:\Users\MMNazari1380\PycharmProjects\Information Retrival\venv\Scripts\python.exe" "C:/Users/MMNazari1380/PycharmProjects/Information Retrival/Phase1
.py"

{'بار': [2, {0: [0, 2]}], 'و': [5, {0: [1, 4, 10, 25, 85]}], 'نغبگ': [1, {0: [6]}], 'دلسوز': [1, {0: [6]}], 'نظا': [2, {0: [6]}]], 'بلکه': [1, {0: [6]]], 'بلکه': [1, {0: [6]]], 'مزار': [7]], 'بلکه': [1, {0: [6]]], 'بلکه': [1, {0: [6]]], 'مزار': [7]], 'بلاه ([13]], 'بله ([14]], 'بله ([14]], 'بله ([14]]], 'بله ([14]]]

Process finished with exit code 0 Oocid

Т