

کار با کتابخانه هضم

گروه ۵

شیوا رادمنش، سحر شیخ الاسلامی، سیدکسری کرمی نژاد، علی نصیری سروی، حامد واعظی

اطلاعات گزارش	چکیده
تاریخ: ۱۳۹۹/۷/۱۹	در این تکلیف به نحوه کار با کتابخانه هضم (hazm) در زبان برنامه نویسی Python پرداخته شد و با استفاده از آن تعدادی سطر مشخص شده را در ابتدا با استفاده از کلاس Normalizer که مربوط به نرمال سازی متون می شود، نرمال سازی کردیم و سپس متن نرمال شده را با استفاده از کلاس های توکن سازی، توکن های آن را جدا نموده و متن را پردازش کردیم. در نهایت برای ۵ جمله تصادفی از متن، کلمات موجود در آن را ریشه یابی کردیم.
واژگان کلیدی:	
توکن سازی	
نرمال سازی	
ریشه یابی کلمات	

۱- مقدمه

بازیابی اطلاعات (IR) پیدا کردن مواد (معمولا اسناد و داکومننت ها) یک گونه غیرساخت یافته (معمولا متن) که یک نیاز اطلاعاتی را از میان مجموعه های بزرگی از داده (که معمولا روی کامپیوترها ذخیره می شوند) تامین می کند. عموما این داده های غیرساخت یافته به فرمت های مختلفی ذخیره شده اند و این می تواند باعث ایجاد مشکل در هنگام بازیابی شود. از این رو لازم است پیش پردازش هایی انجام شود تا اسناد آماده ی استفاده گردند. از انواع این پیش پردازش ها می توان به نرمال سازی، توکن سازی و ریشه یابی کلمات اشاره کرد.

این تکلیف به چهار بخش تقسیم می شود:

در قسمت اول یک فایل با فرمت docx که برای پردازش مناسب نیست را به فرمت txt در می آوریم تا برای پردازش های مراحل بعد مشکلی نداشته باشیم.

در قسمت دوم ابتدا متن مربوط به گروه خود را جدا می کنیم. سپس پیش پردازش های لازم برای متن را اعمال کرده و متن را نرمال سازی می کنیم.

در قسمت سوم ابتدا متن مربوط به گروه خود را به کلمات مجزا شکسته و در نهایت نوع هر کلمه (اسم، فعل و...) را استخراج می کنیم.

در قسمت چهارم ۵ جمله تصادفی از متن مربوط به گروه خود را انتخاب کرده و ریشه کلمات موجود در آن جملات را استخراج می کنیم.

۲- شرح تکنیکال

شرح تکنیکال + کد های گزارش

قسمت اول:

با استفاده از کتابخانه ی "docx2txt" فایل paragraphs.docx را به فایل paragraphs.txt تبدیل کردیم.

قسمت دوم:

در فایلی که در اختیار ما قرار گرفته شده بود، در ابتدای متن هر گروه کلید واژه ی «گروه» به همراه شماره ی گروه نوشته شده بود. برای جدا کردن متن گروه خود (گروه ۵)، متن داده شده را با استفاده از کلید واژه ی «گروه» جداسازی می کنیم.

ماژول "Normalizer" در کتابخانه ی هضم، پیش پردازش هایی را در متن اعمال می کند. این پیش پردازش ها شامل موارد زیر می باشد:

- تبدیل فاصله به نیم فاصله (به وسیله ی تابع "affix_spacing")
- حذف کاف و یای عربی (به وسیله ی تابع "character_refinement")
- ایجاد فاصله ی مناسب بعد (یا قبل) علائم نگارشی (به وسیله ی تابع "punctuation_spacing")

این ماژول شامل تابع "normalize" نیز می باشد که همه ی پیش پردازش های ذکر شده را انجام می دهد. در این پروژه برای اعمال پیش پردازش های لازم از این تابع استفاده شده است.

قسمت سوم:

ابتدای هر پاراگراف، در فایلی که در اختیار ما قرار داده شده بود، کلید واژه‌ی «پاراگراف» به همراه شماره آن پاراگراف داده شده بود. به همین دلیل برای جداسازی هر پاراگراف، ابتدا کل متن گروه خود (گروه ۵) را

بر اساس واژه پاراگراف جداسازی می‌کنیم.

ماژول "[SentenceTokenizer](#)" در این کتابخانه، جملات را شناسایی و جدا می‌کند. با استفاده از این ماژول، جمله‌های هر پاراگراف را استخراج می‌کنیم. با توجه به اینکه بعد از هر پاراگراف در متن، شماره آن پاراگراف آمده است، اولین جمله هر پاراگراف در واقع شماره همان پاراگراف است. به همین دلیل برای بررسی کلمات هر جمله از دومین جمله آن شروع می‌کنیم. خروجی این ماژول، لیستی از جمله‌های هر پاراگراف است و با توجه به توضیحات بالا، یکی کمتر از اندازه‌ی این لیست، تعداد جملات هر پاراگراف را نشان می‌دهد.

ماژول "[WordTokenizer](#)" در کتابخانه‌ی هضم، کلمات هر جمله را شناسایی می‌کند. ورودی این ماژول، جملات می‌باشد و در خروجی کلمه‌های آن جمله در لیست آمده است. با استفاده از این ماژول، کلمات متن را استخراج کردیم. همانند بالا، سائز این لیست تعداد کلمات را نشان می‌هد.

حال با استفاده از ماژول "[POSTagger](#)" نوع هر کلمه را مشخص می‌کنیم. این ماژول یک pattern یا مدل به عنوان ورودی می‌گیرد. ما از مدل "postagger.model" که مدل ارائه شده توسط خود کتابخانه است استفاده کردیم. ورودی این ماژول لیستی از کلمات است و در خروجی، کلمه و نوع آن کلمه (فعل، اسم، ...) را برمی‌گرداند.

پس از مشخص کردن نوع هر کلمه، با شمردن کلماتی که tag آنها برابر "V" می‌باشد، تعداد افعال پاراگراف را بدست می‌آوریم. با انجام همین عمل برای اسم (آنهایی که tag شان برابر "NE" یا "N" باشد)، تعداد اسم‌ها را بدست می‌آوریم.

لیست‌های "verbs_list" و "nouns_list" و "token_list"، به ترتیب شامل کل افعال، اسم‌ها و کلمات متن هستند. همچنین محتوای این لیست‌ها در فایل ذخیره شده است.

خروجی نهایی کد نیز، جدولی است که به ازای هر پاراگراف، تعداد جملات، تعداد کلمات، تعداد افعال و اسم را گزارش کرده است. این جدول ۱ آورده شده است.

Par. No	Sent #	word #	verbs#	nouns#
1	2	77	8	11
2	4	122	9	19
3	3	116	15	22

جدول ۱. تعداد جملات، کلمات، افعال و اسم‌های هر پاراگراف

قسمت چهارم:

در مرحله اول تمامی جملات سه پاراگراف در کنار هم قرار داده تا لیستی از ۹ جمله داشته باشیم. سپس ۵ جمله را به صورت زدنوم از بین آن‌ها انتخاب می‌کنیم.

با استفاده از ماژول "[WordTokenizer](#)" که توضیح آن در مرحله قبل داده شده است، توکن‌های مربوط به هر جمله را استخراج می‌کنیم.

سپس به کمک ماژول "[Lemmatizer](#)" ریشه هر کلمه را استخراج می‌کنیم. این ماژول یک کلمه و نوع آن کلمه را گرفته و به کمک آن ریشه کلمه را پیدا می‌کند (نوع کلمه به عنوان پارامتر اختیاری می‌باشد). سپس هر کلمه و ریشه مرتبط را در یک دیکشنری ذخیره می‌کنیم.

برخی از کلمات و افعالی که توسط این ماژول به درستی ریشه یابی شده اند را با هم بررسی می‌کنیم:

کلمه اصلی	ریشه کلمه	نوع کلمه
نخواهند شد	شد#شو	فعل
بازاری	بازار	اسم
کنید	کرد#کن	فعل
دارد	داشت#دار	فعل
نیست	#هست	فعل
بازیهای	بازی	اسم

جدول ۲. برخی از کلمات و ریشه مرتبط آن‌ها