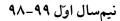
بازيابي ييشرفته اطلاعات





مدرس: حمید بیگی

دانشکدهی مهندسی کامپیوتر

زمان تحویل: ۱ آذر ماه

فاز اول پروژه (۸۵ نمره + ۱۵ نمره امتیازی)

هدف از فاز اول پروژه پیادهسازی یک سیستم بازیابی اطلاعات است. فاز اول پروژه از ۵ بخش تشکیل شده است و دو مجموعه داده نیز در اختیار شما قرار دارد. مجموعه ی اول که به زبان انگلیسی است، شامل چندین مقاله خبری است که توسط موتور جستجوی آکادمیک ComeToMyHead جمعآوری شده اند. توجه کنید که این مجموعه داده شامل دو ستون عنوان و متن خبر است. مجموعه داده دوم به زبان فارسی است و بخشی از پیکره ی ویکی پدیای فارسی شامل چندین صفحه ویکی پدیا به فرمت xml است. توجه نمایید که تمامی صفحات در یک فایل و یکی پدیا و شهمیدن الگوی ذخیرهسازی، هر صفحه را به صورت جدا استخراج نمایید زیرا هر صفحه یک مستند (document) مجزا است.

بخش اول به پیش پردازش متنی داده ها می پردازد که شامل یکسانسازی متن، جداسازی لغات، حذف لغات پرتکرار و ... است. بخش دوم نمایه مسازی است. در بخش سوم نیز باید روی این نمایه فشرده سازی صورت بگیرد. در ادامه، قسمت جستجو و بازیابی سیستم قرار دارد که در بخش چهارم پرسمان ورودی کاربر را باید تصحیح کرده و در بخش پنجم از نمایه های پیاده سازی شده برای جست و جو استفاده می شود.

بخش ۱. پیش پردازش اولیه (۲۵ نمره)

در این بخش از پروژه ابتدا باید مجموعه فایل هایی که در اختیارتان قرار گرفته است را بخوانید سپس به ترتیب مراحل پیش پردازش متنی که در ادامه آمده است را روی آنها اعمال کنید. برای اعمال پیشرو میتوانید از کتابخانههای آماده استفاده کنید. برای زبان پایتون کتابخانه هضم و جاوا JHazm پیشنهاد می شود. برای یکسانسازی متون انگلیسی می توانید از کتابخانه NLTK استفاده کنید.

- ۱. نرمالسازی متنی (normalization): برای یکسان سازی متون میتوانید از توابع کتابخانه های معرفی شده استفاده کنید. اما در صورتی که میخواهید خودتان پیاده سازی کنید باید پیاده سازی تان شامل برگرداندن لغات به ریشه، case folding (برای یکسان سازی متون انگلیسی) و بقیه مواردی که در درس بیان شده است باشد.
 - ۲. جداسازی (tokenization): برای این کار می توانید از توابع کتابخانه های معرفی شده استفاده کنید.
- ۳. حذف علائم نگارشی: هر کدام از مجموعه متنها یک سری علائم نگارشی مثل نقطه، ویرگول و ... دارند
 که آنها را باید حذف کنید.
- ۴. یافتن و حذف لغات پرتکرار (stopwords): در این بخش، حذف درصد معقولی از لغات پرتکرار مورد نظر است. برای این منظور لازم است تا همه متن را پردازش کنید و نسبت به حجم متن، کلماتی که پرتکرار هستند را نمایش دهید. این نسبت را طوری در نظر بگیرید که کلمات پرتکرار به دست آمده، تا حد خوبی منطقی و کافی باشند.
- ۵. بازگرداندن کلمات به ریشه (stemming): در نهایت افعال، اسامی و ... را به حالت ساده و پایه ای خود برگردانید.

بارمبندى

۱. گرفتن متن از کاربر و نمایش لغات آن بعد از پیشپردازش متنی (۱۵ نمره)

۲. نمایش لغات پرتکرار (از متون در اختیار قرار گرفته) (۱۰ نمره)

بخش ۲. نمایهسازی (۲۵ نمره)

در این بخش پیادهسازی نمایه جایگاهی (Positional) و نمایه Bigram مطلوب است. برای نمایه جایگاهی باید به ازای هر لغت، لیستی از اسناد شامل آن لغت و جایگاه(ها) هر لغت در آن سند را داشته باشید و برای نمایه Bigram نیز ترکیبهای دو حرفی تمامی کلمات موجود در لغتنامه که این ترکیب در آنها موجود است را ذخیره کنید. این نمایه برای قسمت اصلاح پرسمان مورد استفاده قرار خواهد گرفت. نمایه شما باید پویا باشد یعنی با حذف سند از نمایه نیز حذف شده و با اضافه کردن سند در طول اجرای برنامه به نمایه اضافه شود. همچنین بعد از نمایهسازی باید قادر باشید نمایه را در فایلی ذخیره کرده و از آن بخوانید.

بارمبندى

- ۱. نمایهسازی از روی پوشههای در اختیار قرار داده شده (۱۵ نمره)
 - ۲. نمایش posting list کلمه ورودی توسط کاربر (۵ نمره)
- ۳. نمایش جایگاه کلمه وارد شده توسط کاربر در هر سند (۵ نمره)

بخش ٣. فشردهسازي نمايهها (١٥ نمره)

در این بخش هدف فشرده سازی نمایه های ساخته شده به دو روش variable byte و gamma code است. (برای ذخیره سازی در فایل و بخش های بعدی می توانید فقط یکی از این دو روش را ادامه دهید.)

بارمبندى

- 1. نمایش میزان حافظه اشغال شده قبل و بعد از اعمال variable byte (۵نمره)
- ٢. نمايش ميزان حافظه اشغال شده قبل و بعد از اعمال gamma code . (۵ نمره)
 - دخیره سازی نمایهها در فایل و بارگذاری از آن (۵ نمره)

بخش ۴. اصلاح پرسمان (۱۰ نمره)

در صورتی که پرسمان ورودی دارای غلط املایی باشد، یا به عبارتی لغت (هایی) از آن در لغت نامه موجود نباشد، لازم است که با جستجوی لغتهای احتمالی و انتخاب بهترین لغت به ادامهی جستجو با پرسمان اصلاح شده پرداخته شود. برای اینکار ابتدا باید به وسیلهی روش bigram و معیار jaccard نزدیک ترین لغات به لغت با غلط املایی را پیدا کنید. سپس بهترین لغت از میان آنها را با استفاده از معیار edit distance بیابید.

بارمبندی

۱. نمایش پرسمان اصلاح شده (۱۰ نمره)

بخش ۵. جستجو و بازیابی اسناد (۱۰ نمره + ۱۵ نمره امتیازی)

در این بخش دو روش جستجو باید به صورتی که در ادامه توضیح داده شده است، پیادهسازی شوند. البته توجه نمایید که روش دوم جستجو امتیازی است.

- ۱. جستجوی ترتیب دار در فضای برداری tf-idf به روش lnc-ltc : در این روش جستجو بعد از دریافت پرسمان ورودی، باید لیستی از اسناد مرتبط به ترتیب امتیاز نمایش داده شود.
- ۲. جستجوی proximity با اندازه ی پنجره ی وارد شده در ورودی (امتیازی): در این روش جستجو ابتدا باید اسنادی که تمام کلمات پرسمان در یک بازه ای به اندازه ی پنجره ی داده شده، در آن سند وجود داشته باشند، پیدا شوند. سپس از بین آن ها به ترتیب امتیازشان براساس جستجوی ترتیب دار در فضای بردار tf-idf به روش Inc-ltc داک ها نمایش داده شوند.

توجه: برای هر دو نوع جستجو نمایش ۱۰ سند در صورت موجود بودن کافی می باشد. مجددا توجه نمایید که روش دوم جستجو امتیازی است.

بارمبندى

- اد. نمایش لیست اسناد مرتبط به ترتیب شباهت در جستجوی ترتیب دار فضای برداری tf-idf به روش ۱ نمایش لیست اسناد مرتبط به ترتیب شباهت در جستجوی ترتیب دار در فضای برداری tf-idf به روش ۱ نمایش لیست اسناد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب دارد نمایش استاد نمایش استاد مرتبط به ترتیب شباهت در جستجوی ترتیب نمایش استاد نما
- ۲. نمایش لیست اسناد مطابق با پرسمان و اندازه پنجره ورودی در جستجوی proximity (۱۵ نمره امتیازی)

بخش ۶. نكات

- 1. باید یک واسط کاربری برای تست موارد مختلف مشخص شده در قسمت بارمبندی هر بخش، برای تحویل حضوری وجود داشته باشد. واسط کاربری می تواند تحت کنسول پیاده سازی شود.
 - ۲. امکان تغییر بارمبندی وجود دارد.