

به نام خدا



دانشگاه تهران

دانشکده فنی

دانشکده مهندسی برق و کامپیوتر



درس بازیابی هوشمند اطلاعات

تمرین دوم

آبان ماه ۱۴۰۲

※ فهرست

بخش ۱ - سوالات عملی	۳
مقدمه	۳
شرح دادگان	۴
پیش‌نیازها	۴
سؤال ۱ - هموارسازی	۵
سؤال ۲ - پیاده‌سازی تابع وزندهی با استفاده از Pseudo Relevance Feedback	۶
بخش ۲ - سوالات تئوری	۸
سؤال ۱ (امتیازی) - الگوریتم Expectation Maximization	۸
سؤال ۲ - مرتب‌سازی اسناد	۹
سؤال ۳ - مرتب‌سازی اسناد با استفاده از هموارساز	۱۰
ملاحظات (حتماً مطالعه شود)	۱۱

بخش ۱- سوالات عملی

مقدمه

در تمرین اول با معیارهای ارزیابی و توابع امتیازدهی به اسناد آشنا شدید. دیدید که یک تابع امتیازدهی با توجه به میزان ارتباط یک سند با پرس‌وجو، امتیازی به سند تخصیص می‌دهد تا در نهایت اسناد براساس امتیازشان، رتبه‌بندی و نمایش داده شوند. در این تمرین قصد داریم روش‌های مختلف هموارسازی توابع بازیابی و پارامترهای آنها را مورد مطالعه قرار بدهیم.

نکات قابل توجه در هنگام پاسخ به سؤالات:

- معیارهای ارزیابی در این تمرین MAP و $P@10$ می‌باشند..
- بدیهی است که حجم تمرین معیار نمره‌ی شما نیست، به تفسیرهایی که بدون آزمایش و صرفاً به صورت فرضی بیان کردند نمره‌ای تعلق نمی‌گیرد.

شرح دادگان

برای این تمرین داده‌های زیر بر روی سایت درس قرار داده شده‌اند.

پیکره متنی^۱ (فایل اسناد):

مجموعه‌ای از ۱۴۰۰ سند به‌دست آمده از چکیده‌های علمی که هر سند شامل فیلدهای زیر می‌باشد:

۱. DOCNO: شناسه هر سند

۲. FILEID: شناسه فایل

۳. HEAD: عنوان سند

۴. TEXT: متن سند

فایل پرس‌وجوها^۲:

این فایل شامل ۱۶۰ پرس‌وجو می‌باشد و فیلدهای زیر را شامل می‌شود:

۱. Number: شناسه پرس‌وجو

۲. Text: متن پرس‌وجو

فایل دادگان طلایی^۳:

این فایل شامل قضاوت‌های مرتبط^۴ می‌باشد در مرحله نهایی جهت ارزیابی کارایی توابع بازیابی مورد استفاده قرار می‌گیرد.

پیش‌نیازها

مشابه تمرین قبلی جهت استفاده از اسناد در توابع بازیابی، بایستی اسناد ابتدا شاخص‌گذاری گردند تا دسترسی به آماره‌های مورد نیاز برای محاسبه‌ی مقادیر امتیازها ساده شود.

هنگام شاخص‌گذاری به نکات زیر توجه کنید:

- نوع فایل را trext قرار دهید.
- از Porter Stemmer جهت ریشه‌یابی کلمات استفاده کنید.
- از Tokenizer جهت جداسازی کلمات موجود در فیلد text استفاده کنید.

¹ Corpus

² Queries

³ Golden Dataset

⁴ Relevance Judgments

سؤال ۱- هموارسازی

یکی از مشکلات حوزه بازیابی اطلاعات و روش‌های Likelihood وجود احتمال‌های صفر برای کلماتی است که در اسناد مشاهده نمی‌شوند. روش‌های هموارسازی برای حل این مشکل مطرح شده‌اند تا احتمال رخداد کلمات دیده نشده پرس‌وجو در اسناد را تخمین بزنند.

ابزار گالاگو، به صورت پیش فرض بازیابی را به روش Query-Likelihood انجام می‌دهد. هدف از این سوال آشنایی با روش‌های هموارسازی می‌باشد. با استفاده از نمودار مناسب برای هر یک از روش‌های خواسته شده مقادیر مختلف λ ، μ و δ را بررسی کنید و مقدار بهینه را گزارش نمایید. (تعداد اسناد بازیابی شده را ۱۰۰ قرار دهید و برای ریشه‌یابی از porter stemmer استفاده کنید).

روش‌هایی که قصد داریم در این تمرین مورد بررسی قرار دهیم عبارتند از:

- روش Additive Smoothing با پارامتر δ
- روش JM با پارامتر λ
- روش Dirichlet Prior با پارامتر μ
- روش هموارسازی دو مرحله‌ای با پارامترهای λ و μ^1

$$p(w|d) = (1 - \lambda) \frac{c(w, d) + \mu p(w|C)}{|d| + \mu} + \lambda p(w|C)$$

راهنمایی:

- با توجه به اینکه روش JM و Dirichlet prior در ابزار گالاگو پیاده‌سازی شده‌است، می‌توانید از توابع پیش‌فرض گالاگو استفاده کنید.
- می‌توانید با تغییر تابع DirichletScoringIterator روش‌ها را پیاده‌سازی کنید.
- در مقداردهی برای پارامترها بهتر است ابتدا گام‌های بلند و سپس گام‌های کوچک آزمایش گردند تا منابع محاسباتی تلف نشود.

^۱ هموارسازی دو مرحله‌ای از ترکیب دو روش هموارسازی Dirichlet و JM بدست می‌آید.

سؤال ۲- پیاده‌سازی تابع وزن‌دهی با استفاده از Pseudo Relevance Feedback

در این سؤال قصد داریم به پیاده‌سازی تابع وزن‌دهی با استفاده از Pseudo Relevance Feedback بپردازیم. برای این منظور فایل‌های زیر در صفحه درس قرار گرفته است:

- wResult.java
- wordWeight.java
- bSearch.java
- fbData.java
- fbMixtureModel.java

فایل‌های فوق شامل پیاده‌سازی روش Mixture Model است، برای پیاده‌سازی این سؤال ابتدا فرمت پرس‌وجوهای موجود در فایل q3.json را به فرمت زیر با پسوند tsv ذخیره کنید:

- #queryNumber [/t(tab)] queryTittle [\n(enter)]

برای مثال:

- 1 what similarity law....

برای انجام این سؤال در ابتدا نیاز است که به API گالاگو متصل شوید. برای این منظور در ابتدا نیاز است که در محیط IDE یک پروژه MAVEN ایجاد گردد. پس از ساخت این پروژه می‌توانید با ویرایش فایل POM.xml و با اضافه کردن API پروژه را به گالاگو متصل کنید. سپس فایل‌های فوق را به پروژه ایجاد شده اضافه نمایید. همچنین نیاز است تا فایل‌های کتابخانه گالاگو موجود در پوشه appassembler را نیز به پروژه خود اضافه کنید.

پس از اتصال به API گالاگو حال می‌توانید با استفاده از تغییراتی بر روی تابع محاسبه وزن مطابق خواسته مسئله، مدل EM را پیاده‌سازی کنید. در نهایت مدل خود را برای پرس‌وجوهای موجود در فایل q3.json اجرا کنید.

الف) با استفاده از نمودار مناسب رابطه بین تعداد سندهای منتخب برای بازخورد به ازای مقادیر بزرگتر از ۱ و معیارهای ارزیابی را نمایش دهید و سپس به تحلیل نتایج بپردازید. در نهایت مقدار بهینه را شناسایی و گزارش کنید.

ب) با توجه به تعداد اسناد منتخب، رابطه بین تعداد کلمات استخراج شده برای بازخورد و معیار ارزیابی را با نمودار مناسب نمایش دهید، مقدار بهینه را گزارش کرده و در نهایت به تحلیل نتایج به‌دست آمده بپردازید.

راهنمایی:

۱. در کلاس `fbMixtureModel.java` تابعی با نام `computeWeights()` ایجاد شده است، می بایست این تابع را به نحوی ویرایش کنید که با استفاده از روش EM، وزندهی به کلمات استخراج شده از سندهای منتخب فراهم شود.
۲. با ویرایش کلاس `bSearch.java` می توانید فایل های ورودی و خروجی را برای پروژه خود تنظیم کنید.

بخش ۲- سؤالات تئوری

سؤال ۱ (امتیازی) - الگوریتم Expectation Maximization

یک مدل مخلوط دو بخشی برنولی^۱ برای تعیین اسناد مرتبط و نامرتبط مانند زیر دارید:

$$P(D|Q) = P(R = 1|Q) P(D|R = 1, Q) + P(R = 0|Q) P(D|R = 0, Q)$$

در این عبارت R نشان دهنده‌ی کلاس ارتباط است. ($R=0$ نامرتبط، $R=1$ مرتبط)

الگوریتم EM در دو گام Expectation و Maximization به تخمین احتمالات $P(R|Q)$ و $P(w|R, Q)$ می‌پردازد و به صورت زیر است:

گام Expectation: مقدار $\gamma(R|D, Q)$ را محاسبه کنید.

گام Maximization: با استفاده از γ پارامترها را دوباره تخمین بزنید.

بعد از ۵ مرحله اجرا^۲ احتمالات کلمه‌ی «retrieval» برابر مقدار زیر است:

$$P(\text{retrieval}|R = 1, Q) = 0.8$$

$$P(\text{retrieval}|R = 0, Q) = 0.3$$

۱. یک سند که شامل ۳ کلمه‌ی retrieval است چگونه دسته بندی می‌شود؟ پاسخ خود را شرح

دهید. (مقادیر $P(R = 1|Q) = P(R = 0|Q) = 0.5$)

راهنمایی: باید مقدار زیر را حساب کنید

$$\gamma(R = 1|D, Q) = P(R = 1|D, Q)$$

سند به صورت:

D: "retrieval retrieval retrieval"

۲. نقطه‌ی ضعف و قدرت این الگوریتم در بازیابی اطلاعات چیست؟

¹ two-component Bernoulli mixture model

² iterations

سؤال ۲- مرتب سازی اسناد

اسناد زیر را در نظر بگیرید

D1: The cat sat on the mat. The cat was black and white.

D2: The dog played with the cat. The animals played in the park.

D3: The girl played in the park with her dog. They played for hours.

شما باید برای پرس و جوهای زیر اسناد را مرتب کنید:

Q1: cat

Q2: cat dog

Q3: cat dog park

۱. Okapi TF^۱ را برای هر کدام از پرس و جوها در هر یک از اسناد محاسبه کنید. مقدار $k1=1.5$

۲. مقدار BM25 را بین هر کوئری و سند با استفاده از پارامترهای زیر محاسبه کنید:

$$k1 = 1.5, b = 0.75, avgdl = 10$$

۳. برای هر کدام از کوئریها و روشها اسناد را بر حسب امتیاز مرتب کنید.

۴. بررسی کنید چه عواملی باعث تفاوت در ترتیب اسناد در این دو روش شده است، همچنین بررسی

کنید چگونه اهمیت کلمات در هر سند مشخص می شود.

¹ Okapi TF/BM25 TF

سؤال ۳- مرتب سازی اسناد با استفاده از هموارساز

اسناد و پرس و جو زیر را در نظر بگیرید:

D1: Machine learning models like regression are used for prediction tasks. Neural networks are a common model.

D2: Models used in machine learning include SVM, decision trees, k-NN. These models have applications in many areas.

Q: machine learning models

۱. مقدار $P(\text{neural}|\text{D1})$ با استفاده از هموارساز JM و مقدار $\lambda=0.3$ به دست آورید.
($P(\text{neural}|\text{C}) = 0.001$)
۲. مقدار $P(\text{SVM}|\text{D2})$ را با استفاده از هموارساز Dirichlet به دست آورید. ($\mu = 1500$)
۳. مقدار query likelihood برای $P(Q|\text{D1})$ و $P(Q|\text{D2})$ با استفاده از هموارساز JM و Dirichlet محاسبه کنید. (مقادیر پارامترها را $\lambda=0.3$ و $\mu = 1500$ در نظر بگیرید)
۴. با استفاده از هموارسازهای متفاوتی که استفاده کردید، برای هر بخش کدام سند بازگردانده می شود؟
۵. برتری های JM را نسبت به Dirichlet برای مرتب سازی اسناد تشریح کنید.

ملاحظات (حتماً مطالعه شود)

- تمامی نتایج شما باید در یک فایل فشرده با عنوان IR_CA2_StudentID تحویل داده شود.
- این فایل فشرده، بایستی حاوی یک فایل با فرمت PDF (گزارش تایپ شده) و یک پوشه به نام Codes باشد که کدهای نوشته شده را به تفکیک هر سؤال شامل شود.
 - خوانایی و دقت بررسی ها در گزارش نهایی از اهمیت ویژه ای برخوردار است. به تمرین هایی که به صورت کاغذی تحویل داده شوند یا به صورت عکس در سایت بارگذاری شوند، ترتیب اثری داده نخواهد شد.
 - در پاسخ به سؤالات عملی، بایستی آزمایش های انجام شده، پارامترهای آزمایش، نتایج و تحلیل ها را به طور کامل شرح دهید.
 - مهلت تحویل تمرین به هیچ عنوان تمدید نخواهد شد. تمرین تا یک هفته بعد از مهلت تعیین شده با جریمه تحویل گرفته می شود که جریمه تأخیر تحویل تمرین تا یک هفته به ازای هر روز ۱۵ درصد است.
 - توجه کنید این تمرین باید به صورت تک نفره انجام شود و پاسخ های ارائه شده باید نتیجه فعالیت فرد نویسنده باشد (همفکری و به اتفاق هم نوشتن تمرین نیز ممنوع است). در صورت مشاهده تقلب به همه افراد مشارکت کننده، نمره تمرین صفر و به استاد نیز گزارش می گردد.
 - در صورت بروز هرگونه مشکل با ایمیل های زیر در ارتباط باشید:

[mailto: mohammad.na3ri@gmail.com](mailto:mohammad.na3ri@gmail.com)

[mailto: mj.kamyab@ut.ac.ir](mailto:mj.kamyab@ut.ac.ir)

مهلت تحویل بدون جریمه: ۱۲ آبان ماه ۱۴۰۲

مهلت تحویل با تأخیر، با جریمه ۱۵ درصد: ۱۹ آبان ماه ۱۴۰۲