

## به نام ایزد منان

تمرین اول درس بازیابی اطلاعات، «روش‌های سنتی بازیابی اطلاعات»



استاد درس: دکتر ممتازی



پاییز ۹۹ - دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر

نکاتی در مورد این تمرین که نیاز به توجه و دقت دوستان دارد:

۱- در جدول زیر نحوه اعمال جریمه تاخیر در ارسال تمرین‌ها ذکر شده است.

میزان جریمه	میزان تاخیر (روز)
هر روز ۵٪	۱ الی ۲ روز
هر روز ۱۰٪	۲ الی ۶ روز

در صورتی که بین ۷ تا ۱۴ روز تاخیر داشته باشید، نمره شما از ۵۰٪ محاسبه می‌شود و پس از این بازه با توجه به سایر تمرین‌ها و زمان تحویل، به تمرین ارسالی نمره‌ای تعلق نمی‌گیرد.

۲- هرگونه کپی کردن باعث عدم تعلق نمره به تمامی افراد مشارکت کننده در آن می‌شود.

۳- آخرین مهلت ارسال تمرین، ساعت ۲۳:۵۵ دقیقه روز چهارشنبه ۱۴ آبان ماه می‌باشد. این زمان با توجه به جمع‌بندی‌های صورت گرفته، شرایط و با توجه به سایر تمرین‌ها در نظر گرفته شده است و **قابل تمدید نمی‌باشد**.

۴- دوستان فایل ارسالی خود را به صورت فشرده و به صورت «شماره دانشجویی\_HW01» مانند HW01\_99131123 نام گذاری کنید. در این فایل باید مواردی نظیر کدها، فایل گزارش و سایر موارد مورد نیاز در هنگام بررسی وجود داشته باشد و صرفاً این فایل در روز ارائه در نظر گرفته می‌شود.

۵- این تمرین دارای **تحویل در محیط اسکایپ** می‌باشد. زمان آن پس از یک هفته از پایان مهلت تمرین از طریق مودل درس اعلام می‌شود.

۶- زبان برنامه‌نویسی این تمرین می‌تواند پایتون، سی‌پلاس‌پلاس و یا جاوا باشد. (پیشنهاد ما پایتون است).

۷- کدهای خود را به صورت مناسب کامنت گذاری کنید. به صورتی که بتوان حداقل روال اجرا و موارد مورد نیاز را درک کرد.

۸- سعی کنید ابتدا تمامی سوالات و بخش‌ها را مطالعه کنید.

۹- استفاده از هیچ کتابخانه آماده‌ای به جز موارد مطرح شده در تمرین مجاز **نمی‌باشد** و شما باید تمامی موارد را پیاده‌سازی کنید.

۱۰- در صورت هرگونه سوال یا مشکل می‌توانید با تدریس‌یار، درس از طریق ایمیل در ارتباط باشید.

## بخش اول – معرفی دادگان

دادگان<sup>۱</sup> ارائه شده در این تمرین شامل ۴ فایل Retrieval, Query, judgement, Corpus می باشد که در ادامه به شرح جزئیات هر کدام می پردازیم.  
دو فایل Retrieval و Corpus از نمونه هایی با فرمت زیر تشکیل شده اند:

DID	شماره یا آیدی نمونه (سند <sup>۲</sup> )
Date	تاریخ ارائه و انتشار سند
Cat	دسته بندی آن سند
Main Text	متن سند

فایل Query شامل کوئری<sup>۳</sup>های مورد استفاده در این تمرین می باشد که هر کوئری یا نمونه این فایل از فرمت زیر تشکیل شده است.

QID	شماره یا آیدی کوئری
Title	عنوان و متن اصلی
Description	توضیحی در مورد کوئری
Narrative	اطلاعات و توضیح مختصر

فایل judgement از سطرهای مانند نمونه زیر و برای ارزیابی و آزمایش مدل ها تشکیل شده است که بیانگر سند مناسب هر کوئری می باشد.

QID	DID	1
-----	-----	---

---

<sup>1</sup> Dataset

<sup>2</sup> Document

<sup>3</sup> Query

### بخش دوم – بازیابی با استفاده از مدل فضای برداری<sup>۴</sup> (۴۵ امتیاز)

در این قسمت شما با استفاده از اسناد موجود در فایل Retrieval اسناد را به صورت TF-IDF نمایش دهید. (برای محاسبه IDF از Corpus که شامل تعداد بیشتری سند می‌باشد، استفاده کنید.) برای این امر دو حالت زیر را در نظر بگیرید.

- ۱- بردار TF اسناد و کوثری‌های به صورت باینری باشد.
  - ۲- بردار TF اسناد و کوثری‌های به صورت شمارشی باشد.
- به طور مثال بردار حاصل برای کوثری «f a a a b c d» به صورت زیر می‌باشد:

نوع بردار	بردار [a, b, c, d, e, f]
باینری	1 1 1 1 0 1
شمارشی	3 1 1 1 0 1

سپس با استفاده از دو معیار فاصله<sup>۵</sup> Cosine و Jaccard ۱۵ سند مشابه به کوثری را گزارش کنید. اگر در کوثری‌ها کلمه‌ای جدید وجود داشت که در دادگان فایل Retrieval وجود نداشت، آن کلمه را در نظر بگیرید. در گزارش خود تحلیل خود را از معیارهای فاصله و ۲ حالت استفاده از بردار TF را شرح دهید.

### بخش سوم – بازیابی با استفاده از مدل احتمالاتی<sup>۶</sup> BM25 (۳۰ امتیاز)

با استفاده از مدل BM25 ۱۵ سند مرتبط برای هر کوثری را بازیابی و گزارش کنید. تاثیر هایپر پارامترها  $b$  و  $k$  را تحلیل کنید. (حداقل سه مقدار برای هر یک از آن‌ها در نظر بگیرید.)

### بخش چهارم – ارزیابی روش‌ها (۲۵ امتیاز)

تمامی روش‌های مطرح شده در بخش دوم و سوم را با استفاده از معیارهای  $P@5$ ,  $P@10$ ,  $MAP^7$  و  $MRR^8$  را ارزیابی و گزارش کنید. تحلیل خود را از نتایج بدست آمده به صورت مختصر بیان کنید.

<sup>4</sup> Vector Space

<sup>5</sup> Distance Metric

<sup>6</sup> Best Match

<sup>7</sup> Mean Average Precision

<sup>8</sup> Mean Reciprocal Rank

## بخش آخر - برخی نکات در مورد گزارش و تمرین

- دادگان مطرح شده در این تمرین و تمامی بخش‌ها را می‌توانید از [لینک](#)<sup>۹</sup> پانوش داندود کنید.
- در این تمرین شما مجاز به استفاده کتابخانه‌های زیر و موارد مشابه و هم‌کاربرد با آن‌ها می‌باشد (تمامی کتابخانه‌های استاندارد پایتون، مجاز می‌باشند):

`numpy, scipy, pandas, genism, pickle, hazm`

- در این تمرین سعی شده است علاوه بر آشنایی شما با کاربرد مباحث ارائه شده در کلاس و لمس بهتر آن، خلاقیت و حل چالش شما نیز ارزیابی شود. لذا در صورتی که در این تمرین چالشی وجود دارد که شما راه حلی برای آن ارائه دادید و استفاده کردید، آن را در گزارش بیان کنید. اما اگر مشکلی بزرگ وجود دارد که نیاز به بررسی مجدد دارد، آن را از طریق ایمیل با تدریس‌یاران درس مطرح کنید.
- در صورتی که هر گونه پیش‌پردازش بر روی دادگان انجام دادید آن را در گزارش خود بیان کنید.
- این تمرین ۱ نمره از بارم کلی شما از تمرینات را با توجه به پوشش مباحث و حجم تمرین دارد. امتیاز این تمرین از ۱۰۰ محاسبه می‌شود که بارم هر بخش مشخص شده است.
- در تمامی بخش‌ها، میزان نتایج شما در ارزیابی شما تاثیر چندانی ندارند (مگر اینکه بسیار دور باشد). بلکه میزان تسلط، دیدگاه و پیاده‌سازی، تحلیل‌ها و خلاقیت شماست که در نمره شما تاثیر مستقیم دارد و بر اساس این موارد مورد ارزیابی قرار می‌گیرد.

موفق باشید - رستمی

<sup>۹</sup> <https://drive.google.com/drive/folders/1cX-6tLr076TtdnhcPCOYlhCRsiyEVCq?usp=sharing>