



دانشگاه فردوسی مشهد  
دانشکده مهندسی  
گروه مهندسی کامپیوتر



تمرین شماره 1  
مبانی بازیابی اطلاعات و جست و  
جوی وب

## پیش‌پردازش متن و نمایه‌گذاری

زمستان 1400، بهار 1401

در این تمرین از شما خواسته شده است تا مراحل پیش‌پردازش و نمایه‌گذاری معکوس<sup>1</sup> را بر روی دیتاستی که به شما معرفی می‌شود، اعمال کنید.

### دیتاست

دیتاست در نظر گرفته شده از این [لینک](#) قابل دسترسی است. این دیتاست، یک فایل با پسوند XLSX می‌باشد که حاوی فراداده‌های کتب یکی از سایت‌های دانلود کتاب آنلاین فارسی است. ویژگی‌های موجود در این فایل به شرح زیر هستند:

---

<sup>1</sup> Inverted indexing

<sup>2</sup> Metadata

- عنوان (title): نام کتاب
- تاریخ (date): تاریخ درج کتاب در سایت
- محتوا (content): خلاصه محتوای کتاب
- دسته‌بندی (category)
- نویسنده (author)
- دیدگاه‌ها (comments)

### پیش‌پردازش

در فاز پیش‌پردازش باید روش‌هایی مانند normalization حذف کلمات توقف<sup>۳</sup>، ریشه‌یابی<sup>۴</sup>، بن‌یابی<sup>۵</sup> و ... (که از دید شما منطقی و لازم است) را متناسب با نیازهای مسئله بر روی ویژگی‌(ستون) «محتوا» در دیتاست اعمال کنید. در پایان فاز پیش‌پردازش باید توکن‌های «محتوا» مربوط به هر یک کتب را استخراج کرده باشید.

### نمایه‌گذاری معکوس

در این فاز، باید توکن‌های استخراج شده از فاز پیش‌پردازش را به صورت معکوس نمایه‌گذاری کنید. برای نمایه‌گذاری کردن می‌توانید از شناسه‌های افزایشی<sup>۶</sup> استفاده کرده و یا با تولید اعداد تصادفی و منحصر به فرد، شناسه تولید کنید.

در پایان فاز نمایه‌گذاری معکوس، پیاده‌سازی شما باید بتواند یک واژه<sup>۷</sup> را به عنوان ورودی دریافت کرده و لیست شناسه کتاب‌هایی که آن واژه را در محتوای خود دارند، خروجی دهد.

**نکته:** برای پیاده‌سازی، محدودیتی در استفاده از زبان یا کتابخانه وجود ندارد.

---

<sup>3</sup> Stop words

<sup>4</sup> Stemming

<sup>5</sup> Lemmatization

<sup>6</sup> Incremental

<sup>7</sup> Term

## موارد تحویلی

فایل ارسالی شما، یک فایل فشرده شده (zip یا rar) خواهد بود که باید شامل موارد زیر باشد:

1- گزارش و مستندات کد پیاده سازی (pdf یا word)

توضیح مراحل پیاده سازی و قطعه کدهای مربوط به هر قسمت (ترجیحا عکس کد خود را به

گزارش اضافه کنید)، ذکر مشخصات محیط پیاده سازی و ماژول های اصلی، بیان دلایل استفاده

از روش ها و ابزارها، بیان چالش ها، تعیین فعالیت های انجام شده توسط هر یک از اعضای گروه و

...

2- کد پیاده سازی

3- فایل خروجی فاز نمایه گذاری معکوس (csv یا xls)

در نهایت فایل فشرده شده را به صورت HW01-GroupX نام گذاری کنید. به جای X، شماره گروه خود را

جایگزین نمایید.

موفق باشید.