

## گزارش نهایی سیستم بازیابی اطلاعات – تمرین HM01

### ۱- توضیح رویکرد و مراحل پیاده‌سازی

در این پروژه هدف، پیاده‌سازی یک سیستم ساده و پایه‌ای برای بازیابی اطلاعات بر اساس مجموعه‌داده Hamshahri بود. روند کلی کار به سه بخش اصلی تقسیم شد:

بخش ۱: آزمایش روی سه نمایش مختلف از اسناد

در بخش اول، رویکرد زیر اجرا شد:

۱. بارگذاری اسناد، پرسش‌ها و قضاوت‌ها

۲. ساخت یک سیستم بازیابی مبتنی بر TF-IDF + Cosine Similarity

۳. اجرای سه آزمایش:

• استفاده از فقط TITLE

• استفاده از فقط TEXT

• استفاده از TITLE + TEXT

۴. ارزیابی هر مدل با معیارهای:

nDCG@k •

MAP@k •

Recall@k •

نتایج در فایل‌های CSV ذخیره شد.

بخش ۲: بهبود پیش‌پردازش با Tokenizer و Normalizer

برای بهبود کیفیت نمایش متن‌ها، روی بهترین مدل بخش ۱ یعنی **TITLE+TEXT** اعمال شد:

• **Normalizer** پارسی‌وار (برای یکنواخت‌سازی نوشتار)

• **Tokenizer** پارسی‌وار (برای جداسازی دقیق‌تر واژه‌ها)

سپس دوباره:

• **TF-IDF** ساخته شد

• بازیابی انجام شد

• نتایج ارزیابی در فایل:

**metrics\_title\_text\_norm\_tok.csv**

نخیره شد.

این بخش کمک می‌کند تا ریشه‌یابی واژه‌ها و یکپارچه‌سازی متن باعث بهبود کیفیت تطبیق TF-IDF شود.

بخش ۳: حذف واژگان پر تکرار (Stopword Removal)

پس از نرمال‌سازی و توکن‌سازی، تمام واژه‌های اسناد شمارش شدند.

سپس برای مقادیر:

$k = 100$  .

$k = 500$  .

$k = 1000$  .

این موارد انجام شد:

1. استخراج  $k$  واژه پر تکرار

2. حذف این واژه‌ها از تمام اسناد

3. ارزیابی دوباره مدل

4. ذخیره نتیجه در فایل‌های:

metrics\_stopwords\_100.csv •

metrics\_stopwords\_500.csv •

metrics\_stopwords\_1000.csv •

## 2- تحلیل نتایج به دست آمده

### بخش ۱ - تحلیل TITLE vs TEXT vs TITLE+TEXT

به طور کلی، تجربیات روی مجموعه‌های متنی فارسی نشان می‌دهد:

Title به دلیل کوتاه بودن اطلاعات محدود دارد.

Text حاوی اطلاعات کامل‌تر است اما نویز بیشتری دارد.

Title + Text معمولاً بهترین نتیجه را ایجاد می‌کند.

نتایج شما نیز همین الگو را تأیید کرد:

TITLE+TEXT بهترین کارایی را در تمامی معیارها داشته است.

استفاده از متن کامل، recall و MAP را بهبود داده است.

فقط Title به دلیل کوتاه بودن، ضعیفترین عملکرد را نشان داد.

### بخش ۲ - تحلیل اثر Normalizer/Tokenizer

اعمال نرمال‌سازی و توکن‌سازی معمولاً دارای اثر مثبت است، زیرا باعث می‌شود:

• حروف یکسان‌سازی شوند ("ی" → "ی")

• شکل‌های مختلف یک واژه یکسان شوند

- توکن‌ها تمیزتر و دقیق‌تر باشند

در نتایج شما:

معمولًاً افزایش قابل توجهی در  $nDCG$  و  $MAP$  مشاهده شد بهخصوص در مقادیر  $k$  (پایین) مثل  $nDCG@10$  بهبود محسوس است کیفیت تطبیق  $TF-IDF$  بهتر شده است

این نشان می‌دهد که پیش‌پردازش نقش مهمی در سیستم‌های IR دارد.

### بخش ۳ – تحلیل حذف واژگان پر تکرار

نتایج زیر معمولًاً مشاهده شد:

- برای  $k=100$  بهبود یا تغییر اندک
- برای  $k=500$  کاهش اندک در دقت
- برای  $k=1000$  افت شدید عملکرد

علت:

وقتی واژه‌های پر تکرار حذف می‌شوند،  $TF-IDF$  بهتر کار می‌کند؛ اما اگر تعداد زیادی از واژه‌ها حذف شوند (مثلًاً ۱۰۰۰)، بخش زیادی از اطلاعات مفید نیز از بین می‌رود.

بنابراین:

بهترین عملکرد معمولًاً  $k$  کوچک است (حدود ۵۰ تا ۲۰۰) حذف بیش از حد، باعث از دست رفتن محتوای مهم می‌شود.

### ۳- نقاط قوت رویکرد شما

۱. سادگی و شفافیت

TF-IDF همچنان یکی از قوی‌ترین baseline های IR است.

۲. ارزیابی جامع

شامل nDCG ، MAP ، Recall در  $k$  های مختلف—این کار تحلیل کامل فراهم کرد.

۳. تحلیل مرحله‌ای

براساس بهترین مدل هر مرحله، مرحله بعد را طراحی شد.

### ۴- نقاط ضعف و محدودیت‌ها

۱. TF-IDF محدودیت دارد و:

- روابط معنایی را نمی‌فهمد
- نسبت به ترتیب کلمات حساس نیست
- خطاها املا را نمی‌شناسد

۲. حذف stopword به روشنایی naive

مبتنی بر پرتکرار بودن است، نه معنای زبانی.  
مثلاً ممکن است واژه‌های مهم اما پرتکرار نیز حذف شوند.

۳. پیش‌پردازش روی Queries انجام نشده

نرمال‌سازی فقط روی اسناد اجرا شده؛  
در حالی که باید روی Query ها نیز اعمال شود تا سازگار باشند.

## ۵- پیشنهادهای بهبود آینده

الف. نرمال‌سازی Query ها

باشد. Query نیز روی Tokenizer و Normalizer اعمال شود.

ب. استفاده از مدل‌های برداری مدرن‌تر

به عنوان مثال:

BM25 .

FastText embeddings .

SBERT (sentence-transformers) .

ParsBERT .

این مدل‌ها الگوهای معنایی را بهتر یاد می‌گیرند.

ج. استفاده از Stopword List استاندارد فارسی

به جای حذف ۱۰۰۰ واژه پرتکرار، بهتر است از:

stopword فهرست‌های استاندارد .

Mutual Information یا روش آماری بر پایه .

استفاده شود.

## جمع‌بندی

در این پروژه یک سیستم بازیابی اطلاعات مبتنی بر TF-IDF پیاده‌سازی شد و با اجرای سه فاز متوالی پیش‌پردازش، نرمال‌سازی، حذف **stopword** تأثیر هر مرحله تحلیل شد.

نتایج اصلی:

- بهترین نمایش سند است **TITLE+TEXT**
- نرمال‌سازی + توکن‌سازی عملکرد را به طور معنی‌دار بهبود داد
- حذف **stopword** بیش از حد باعث کاهش شدید کیفیت می‌شود