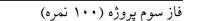
# بازيابي ييشرفته اطلاعات

## نيمسال اوّل ٠٠ ـ ٩٩

زمان تحویل: ۱۵ بهمن ماه

دانشكدهى مهندسي كامپيوتر مدرس: حميد بيگي



هدف از فاز سوم پروژه پیادهسازی الگوریتمهای خوشهبندی و یک خزنده برای واکشی اطلاعات مقالات از سایت Microsoft Academic است. در بخش اول این فاز الگوریتمهای خوشهبندی را برای یک مجموعه دادهی فارسی پیادهسازی میکنید و خوشههای به دست آمده را در خروجی برمیگردانید. در بخش دوم نیز یک خزنده برای واکشی اطلاعات مقالات از سایت Microsoft Academic پیادهسازی میکنید و در آخرین بخش PageRank را برای مقالات واکشی شده محاسبه میکنید.

## بخش ۱. خوشهبندی (۵۰ نمره)

در بخش اول باید برای یک مجموعه داده به زبان فارسی، چند الگوریتم خوشهبندی را پیادهسازی کنید. مجموعه داده انتخابی شامل اطلاعات یک مجموعه اخبار در قالب یک فایل json است. در این دادگان، اطلاعات زیر برای هر خبر موجود است:

- title: عنوان خبر
- summary: خلاصه متن خبر
- link: لینک به خبر در وبسایت
- tags: دستهبندی خبر. به فرمت «دسته اصلی > دسته فرعی». مثلا: «دانش > پزشکی»

برای تبدیل متون اخبار به فضای برداری باید یک بار از طریق TF-IDF و یک بار از Word2vec استفاده کنید. توجه کنید که برای هر دو روش میتوانید از توابع و کتابخانههای آماده استفاده کنید. ضمنا استفاده از کتابخانه هضم برای پیش پردازش متون توصیه می شود. الگوریتمهای خوشه بندی که باید پیاده سازی شوند نیز عبارت هستند از:

- K-means .\
- Gaussian Mixture Model . Y
  - Hierarchical clustering . T

توجه: انتخاب تمامی پارامترهای الگوریتمهای بالا برعهده ی خودتان است. نحوه انتخاب بهترین عدد برای تعداد خوشهها و مقادیر سایر پارامترهایی که انتخاب میکنید را در گزارش بیاورید. برای پیادهسازی الگوریتمهای خوشه بندی نیز می توانید از توابع و کتاب خانههای آماده استفاده کنید. به ازای هر زوج از روشهای تبدیل به فضای برداری و الگوریتم خوشه بندی یک فایل csv در خروجی داشته باشید (مجموعا ۶ فایل) که نتیجه ی خوشه بندی الگوریتم شما است. این فایل ها باید دو ستون داشته باشند: ستون اول لینک خبر و ستون دوم شماره ی خوشه برای آن خبر را قرار دهید. همچنین در گزارش خود تعدادی نمودار از نتایج خوشه بندی های خود (در قالب نمودار دوبعدی یا دندروگرام و یا هر نمودار جالب دیگری) ارائه کرده و در چند سطر به صورت مختصر مشاهدات خود به همراه توضیحاتی ارائه دهید. همچنین حداقل ۲ معیار ارزیابی خوشه بندی را به ازای هر بار خوشه بندی، در گزارش بیاورید. در حین ارزیابی، از دسته بندی اصلی اخبار به عنوان ground truth استفاده کنید.



#### بارمبندى

بارمبندی این بخش قطعی نیست. بر حسب امتیازهایی که خوشهبندیهای شما به دست میدهند نمره تان تعیین میگردد.

## بخش ۲. پیادهسازی خزنده، واکشی اطلاعات مقالات (۴۰ نمره)

در این بخش قصد داریم تا برای سایت Microsoft Academic یک خزنده پیادهسازی کرده و با استفاده از آن اطلاعات تعدادی مقاله را واکشی کنیم.

اطلاعاتی که از هر مقاله باید جمع آوری شوند عبارت هستند از:

- ١. عنوان مقاله
- ۲. چکیدهی مقاله
- ٣. سال انتشار مقاله
- ۴. تمامى نويسندگان مقاله
- ۵. ارجاعات مقاله. توجه کنید که تنها ۱۰ ارجاع اول که در صفحهی مقاله در سایت Microsoft Academic ه. قرار دارد کافی است و نیازی به واکشی تمامی ارجاعات نیست.

خزنده برای آغاز کار باید از چند مقالهای که در فایل start.txt وجود دارند و در صف خزش قرار می گیرند شروع کرده و ۵۰۰۰ مقاله (تعداد کل مقالات به عنوان پارامتر ورودی داده می شود) را ذخیره نماید. همچنین آدرس ۱۰ مقالهی ابتدایی در لیست ارجاعات مقالهی کنونی به صف خزش خزنده اضافه می شود. توجه نمایید که برخی مقالات دارای چند نسخه منتشر شده در وبسایت های مختلف هستند و یکی از این نسخه ها باید ذخیره شود. هیچ مقاله این نبید بیش از یک بار ذخیره شود. اگر لیست ارجاعات یک مقاله کمتر از ۱۰ مورد باشد ایرادی ندارد. همچنین برخی از ارجاعات به صورت لینک نیستند که می توانید از آن ها چشم پوشی کنید.

یک نمونه فایل json از اطلاعات ذخیره شده در فایل sample.json قرار دارد. برای ذخیرهسازی اطلاعات مقالهها مشابه این فایل نمونه اقدام نمایید. توجه کنید که برای خزش سایت Microsoft Academic شاید نیاز باشد تا بین درخواستهای خود تاخیر (delay) بیاندازید.

#### بارمبندى

- ۱. پیادهسازی خزنده (۳۰ نمره)
- ۲. ذخیرهی اطلاعات مقالههای به فرمت json و بررسی اولیه اطلاعات استخراج شده (۲۰ نمره)

# بخش ۲۰ PageRank (۱۰ نمره)

در آخرین بخش الگوریتم PageRank را بر روی مقالات واکشی شده اجرا کرده و نتایج آن را به دست می آوریم. PageRank را به دست می آوریم. PageRank به اگر مقاله PageRank به مقاله PageRank ارجاع (reference) داشته باشد، آنگاه پیوندی PageRank از مقاله PageRank مورد نیاز در ورودی گرفته می شود و سپس معیار PageRank می مقاله PageRank در نظر می گیریم. برای این منظور مقدار PageRank می می مقالات محاسبه شده و در خروجی چاپ می شود. برای این بخش می توانید از ابزارهای آماده استفاده کنید.

### بارمبندى

۱. محاسبهی معیار PageRank برای مقالات واکشی شده و گزارش مقالات با بالاترین رنک (۱۰ نمره)

# بخش ۴. نکات

- ۱. امکان تغییر بارمبندی وجود دارد.
- ۲. نوشتن گزارش فراموش نشود. به قوانین کلاس و پروژه که در پیاتزا قرار گرفته است رجوع کنید.