



## کاوش داده گان انبوه

دستیاران آموزشی

مهدی صادقی (mahdisadeghi@ut.ac.ir)

شیوا پارسا راد (shiva.parsarad@ut.ac.ir)

علیرضا جعفری (Alirezajafari@ut.ac.ir)

ایمان برازنده (barazandeh.iman@ut.ac.ir)

دکتر سامان هراتی زاده

دانشگاه تهران - دانشکده علوم و فنون

ترم اول سال تحصیلی - پاییز ۱۳۹۹

## پروژه نهایی: این ترم هم به آخر رسید

مهلت تحویل این تمرین سه شنبه ۵ اسفند ساعت ۲۳:۵۹ می باشد.

این ترم هم بالاخره تمام شد. به عنوان پروژه نهایی می توانید از بین دو تمرین زیر یکی را انتخاب کرده و پیاده سازی کنید.

**(تمرین اول):** موتور جستجوی تمرین چهارم را (روی همان دیتا ست) با استفاده از Map-Reduce پیاده سازی و روی پلتفرم هادوپ اجرا کنید. برای ارتباط با هادوپ از پایتون، می توانید از کتابخانه MRJob استفاده کنید.

**(تمرین دوم: تا ۲۵٪ نمره اضافی):** می خواهیم یک سیستم توصیه گر مبتنی بر گراف بنویسیم. دیتاست این تمرین نیز MovieLens100k است. (تعداد کاربران را محدود به کارایی سخت افزار کامپیوتر خود کنید). کتابخانه Networkx پایتون کار شما را برای پیاده سازی راحت می کند.

برای پیاده سازی این سیستم بصورت زیر عمل کنید:

- یک گراف دو بخشی<sup>۱</sup> کاربر-آیتم بسازید بطوریکه یک بخش آن را کاربران و یک بخش آن را آیتم ها تشکیل می دهند. هر کاربر به شرطی به یک آیتم لینک می شود که امتیازش به آن بیشتر یا مساوی ۳ (با مقادیر متفاوت تست کنید و بهترین را انتخاب کنید) باشد.
- گراف را Fold کرده و آن را تبدیل به یک گراف کاربر-کاربر کنید. کاربرانی به هم متصل می شوند که حداقل به ۳ محصول بطور مشترک امتیاز داده باشند. (با مقادیر متفاوت تست کنید و بهترین را انتخاب کنید)
- با استفاده از یکی از الگوریتم های کلاسترینگ (اجتماع یابی<sup>۲</sup>) مثل Louvain کاربران مشابه را پیدا کنید.
- آیتم هایی که کاربران در هر کلاستر به آنها امتیاز داده اند در یک لیست جمع کنید.

<sup>1</sup> Bipartite

<sup>2</sup> Community Detection



## کاوش داده‌گان انبوه

دستیاران آموزشی

مهدی صادقی (mahdisadeghi@ut.ac.ir)

شیوا پارسا راد (shiva.parsarad@ut.ac.ir)

علیرضا جعفری (Alirezajafari@ut.ac.ir)

ایمان برازنده (barazandeh.iman@ut.ac.ir)

دکتر سامان هراتی زاده

دانشگاه تهران - دانشکده علوم و فنون

ترم اول سال تحصیلی - پاییز ۱۳۹۹

- برای هر کاربر آیتم‌هایی که مشاهده کرده را از آن حذف کنید و آیتم‌های باقیمانده را بر اساس معیار خاصی مرتب کنید (هر معیاری! مهم نیست، مثلاً براساس مجموع امتیازهایی که این آیتم‌ها دارند یا تعداد امتیازها، هر چه معنی‌دارتر بهتر).
  - به هر کاربر Top-5 آیتم از این لیست را پیشنهاد دهید.
- گراف‌ها در هر مرحله و همچنین خوشه‌ها بصورت گرافیکی نیز نمایش داده شوند. لیست پیشنهادها را نیز برای چند کاربر گزارش کنید.

## نکات تکمیلی

- از زبان برنامه‌سازی پایتون کنید.
- کدها با نامگذاری مناسب در پوشه SOURCES قرار گیرد.
- گزارش کامل تمرین را با ساختار مناسب، بدون ابهام و به زبان ساده نوشته و با فرمت PDF در پوشه دیگری قرار دهید. برای نگارش گزارش‌ها نیز امتیاز مثبت و منفی در نظر گرفته شده است، پس برای نگارش گزارش خود وقت مناسبی اختصاص دهید.
- پاسخ تمرین خود را به ایمیل [barazandeh.iman@ut.ac.ir](mailto:barazandeh.iman@ut.ac.ir) ارسال نمایید. هر دو پوشه را در یک فایل ZIP با نام MMDS\_STID\_FullName\_FP مثلاً MMDS\_830123451\_ImanBarazandeh\_FP قرار داده و در ایمیلی با همین عنوان نیز ارسال کنید.
- در صورت مشخص شدن هر نوع تقلب انفرادی یا مشترک نمره تمرین صفر در نظر گرفته می‌شود.
- به ازای هر روز تاخیر ۲۵٪ از امتیاز تمرین از دست خواهید داد.
- هر گونه ابهام را با دستیاران آموزشی مطرح کنید.