



## کاوش داده‌گان انبوه

دستیاران آموزشی  
مهدی صادقی (mahdisadeghi@ut.ac.ir)  
شیوا پارسا راد (shiva.parsarad@ut.ac.ir)  
ایمان برازنده (barazandeh.iman@ut.ac.ir)

دکتر سامان هراتی زاده  
دانشگاه تهران - دانشکده علوم و فنون  
ترم اول سال تحصیلی - زمستان ۱۳۹۹

### مجموعه تمرین پنجم :

## شناسایی مجموعه اقلام پرتکرار و گاهی نیز کم تکرار!

مهلت تحویل این تمرین سه شنبه ۳۰ دی ساعت ۲۳:۵۹ می‌باشد.

هدف از این تمرین آشنایی با یکی از تکنیک‌های بنیادی داده کاوی یعنی کاوش مجموعه اقلام پرتکرار است. کاوش مجموعه اقلام پرتکرار یکی از تکنیک‌های اکتشافی بدون ناظر است که برای شناسایی همبستگی‌های با اهمیت بین اقلام درون داده‌گان انبوه استفاده می‌شود. این تکنیک این روزها کاربردهای جذابی در بیوانفورماتیک، امنیت و حریم خصوصی، ساختن سبد سهام<sup>۱</sup> و زمینه‌های بسیار دیگر دارد.

### ۱- (طراحی و پیاده سازی الگوریتم)

داده مورد استفاده در این تمرین شامل سبدهای خرید یک خرده فروشی است. هر سطر از این دیتاست نشانگر یک سبد خرید و در هر سطر تعدادی شناسه محصول وجود دارد که نشانگر محصولات خریداری شده است. در این تمرین شما فقط با `product_id` و `order_id` کار خواهید داشت.

الف) با استفاده از الگوریتم Apriori و  $\text{support}=0.03$  کالاهای پرتکرار در سبدهای خرید را پیدا کنید.

ب) با استفاده از الگوریتم PCY و  $\text{support}=0.04$  کالاهای پرتکرار در سبدهای خرید را پیدا کنید.

---

<sup>1</sup> Stock Portfolio



## کاوش داده‌گان انبوه

دستیاران آموزشی

مهدی صادقی (mahdisadeghi@ut.ac.ir)

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

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

دکتر سامان هراتی زاده

دانشگاه تهران - دانشکده علوم و فنون

ترم اول سال تحصیلی - زمستان ۱۳۹۹

## ۲- (پژوهش)

**الف)** یکی از بهبودهایی که روی تعریف کلاسیک مجموعه اقلام پرتکرار برای بعضی کاربردها ضروری است، اضافه شدن وزن به اقلام درون مجموعه‌ها است، به عبارت دیگر مجموعه اقلام پرتکرار وزن دار<sup>۲</sup>.

چرا وزن؟ بدیهی است که در بسیاری از کاربردهای دنیای واقعی، اقلام درون مجموعه‌ها اهمیت و ارزش یکسانی ندارند. بعنوان مثال، سهام شرکت‌های مختلف درون یه یک سبد سهام را در نظر بگیرید که هر کدام نرخ سود متفاوتی در یک روز معاملاتی دارند. با بررسی یک تحقیق کاربردی معتبر، الگوریتم شناسایی مجموعه اقلام پرتکرار وزن دار و همچنین کاربردی که از آن استفاده کرده را با جزئیات کافی به ما نیز معرفی کنید.

**ب)** شاید برایتان جالب باشد که یک زمینه کاربردی و تحقیقاتی به نام **کاوش مجموعه اقلام کم تکرار<sup>۳</sup>!!!!** نیز وجود دارد. کمی در این مورد تحقیق کرده و الگوریتم‌ها و کاربردهای آن را گزارش کنید.

<sup>۲</sup> Weighted Frequent Itemset

<sup>۳</sup> Infrequent Itemset Mining



## کاوش داده‌گان انبوه

دستیاران آموزشی  
مهدی صادقی (mahdisadeghi@ut.ac.ir)  
شیوا پارسا راد (shiva.parsarad@ut.ac.ir)  
ایمان برازنده (barazandeh.iman@ut.ac.ir)

دکتر سامان هراتی زاده  
دانشگاه تهران - دانشکده علوم و فنون  
ترم اول سال تحصیلی - زمستان ۱۳۹۹

## نکات تکمیلی

- از زبان برنامه سازی 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\_HW# مثلا MMDS\_830123451\_ImanBarazandeh\_HW1 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.