

به نام خدا



دانشگاه صنعتی امیرکبیر  
( پلی تکنیک تهران )

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

## تمرین سری سوم داده کاوی

### توضیحات:

- پاسخ به تمرین‌ها باید به صورت انفرادی صورت گیرد و در صورت مشاهده هرگونه تقلب نمره صفر برای کل تمرین منظور خواهد شد.
- تمیزی و خوانایی گزارش تمرین از اهمیت بالایی برخوردار است.
- گزارش تمرین خود را در قالب یک فایل PDF با نام «**HW3\_StudentNumber.pdf**» به همراه کد های بخش پیاده سازی (فایل های `ipynb` یا `.py` ) در فایلی به نام «**HW3\_StudentNumber.zip**» قرار داده و در سایت درس در مهلت معین بارگزاری نمایید.
- توجه داشته باشید که به سوالات پیاده سازی بدون گزارش نمره ای تعلق نمی گیرد.
- در صورت داشتن اشکال می‌توانید از طریق ایمیل **datamining.fall2020@gmail.com** با تدریس‌یاران درس در ارتباط باشید.
- همچنین لازم بذکر است که اگر مواردی در کلاس تدریس نشده انتظار می رود که خود دانشجویان جستجو کنند و انجام دهند.

سوال ۱- مزایا و معایب روش های دسته بندی مشتاق<sup>1</sup> (مثل درخت تصمیم، بیزین، شبکه عصبی) و روش ها دسته بندی تنبل<sup>2</sup> (مثل KNN) را مقایسه کنید.

سوال ۲-

الف) تراکنش های زیر را در نظر بگیرید. با فرض  $\min\_sup = 60\%$  ,  $\min\_conf = 80\%$  تمام frequent itemset ها را با استفاده از الگوریتم های Apriori و FP-growth پیدا کنید.

TID	Items_bought
T1	{M, O, N, K, E, Y}
T2	{D, O, N, K, E, Y}
T3	{M, A, K, E}
T4	{M, U, C, K, Y}
T5	{C, O, O, K, I, E}

ب) الگوریتم های Apriori و FP-growth را از لحاظ بهینگی عملکرد مقایسه کنید.

سوال ۳- در بسیاری از موارد در دنیای واقعی، برای انجام عمل دسته بندی، بخش زیادی از داده های موجود در پایگاه های داده برچسب مشخصی ندارند. از جمله روش های پیشنهادی برای حل این مشکل، روش های مبتنی بر یادگیری نیمه نظارتی<sup>3</sup>، یادگیری فعالانه<sup>4</sup> و یادگیری انتقالی<sup>5</sup> می باشند. درباره چگونگی این روش ها، موارد استفاده هر کدام و همچنین چالش های آن ها به اختصار توضیح دهید.

---

<sup>1</sup>eager

<sup>2</sup>lazy

<sup>3</sup>Semi-supervised learning

<sup>4</sup>Active learning

<sup>5</sup>Transfer learning

**سوال ۴-** دو دسته اصلی روش های یادگیری گروهی<sup>۶</sup> را نام برده و تفاوت آن ها را توضیح دهید. مشخص نمایید هر یک از روش های پیشنهادی در بهبود چه مشکلی حین آموزش مدل یادگیری موثر هستند؟

**سوال ۵-** در ارتباط با ماشین های بردار پشتیبان<sup>۷</sup> به سوالات زیر پاسخ دهید.

الف) توجه به SVM دو کلاسه شرح داده شده در کلاس، یک سناریو برای SVM چند کلاسه (M کلاس) با ذکر توضیحات ارائه دهید.

ب) منظور از مدل با حاشیه سخت<sup>۸</sup> چیست؟

ج) فرض کنید یک ماشین بردار پشتیبان با مرزبندی خطی آموزش داده اید و متوجه می شوید دچار کم برازش شده است. برای حل این مشکل پارامترهای مدل خود را چگونه تغییر می دهید؟

**سوال ۶-** جدول زیر نوع کتاب های خریده شده یک فروشگاه را نشان می دهد:

نوع کتاب خریداری شده	شماره خرید
مذهبی، رمان، شعر	۱
مذهبی، رمان، تاریخی	۲
روانشناسی، رمان، تاریخی، شعر	۳
روانشناسی، تاریخی، شعر	۴
روانشناسی، مذهبی، رمان	۵
روانشناسی، مذهبی، رمان، تاریخی	۶

<sup>۶</sup>Ensemble

<sup>۷</sup>Support Vector Machines

<sup>۸</sup>Hard Margin

۷	روانشناسی
۸	روانشناسی، مذهبی، رمان
۹	روانشناسی، مذهبی، تاریخی
۱۰	روانشناسی، مذهبی

الف) با استفاده از روش FP-Growth تمامی frequent itemset هایی که به کتاب نوع رمان ختم می شود را بیابید. (support = ۲۰٪).

ب) با در نظر گرفتن confidence = ۵۰٪ قواعد معتبر قابل استخراج از frequent itemset های به دست آمده از قسمت الف را بیابید.

## پیاده سازی:

### سوال ۸:

در این بخش هدف استفاده از جنگل تصادفی برای کلاس بندی می باشد. مسئله تایتانیک (سوال ۶ در تمرین ۲) را این بار با به جای درخت تصمیم با استفاده از جنگل تصادفی پیاده سازی کرده و دقت پیش بینی را اندازه گیری کنید (۲ حالت برای عمق درخت ها و ۲ حالت برای تابع تقسیم gini یا entropy در مجموع ۴ حالت مختلف). سپس مقایسه های زیر را انجام دهید:

۱- بالاترین دقت به دست آمده توسط مدل جنگل تصادفی را با مدل درخت تصمیم در تمرین قبل مقایسه کنید.

۲- سرعت یادگیری و تست جنگل تصادفی و درخت تصمیم را با هم مقایسه کنید.

**توجه ۱-** همانطور که در تمرین قبل اطلاع رسانی شد، برای محاسبه دقت پیش بینی روی مجموعه داده تایتانیک یکی از دو روش زیر را می توانید انتخاب کنید:

۱. تقسیم بندی داده های فایل train.csv به دو قسمت train set و test set (با نسبت ۸۰ به ۲۰) و

محاسبه دقت مدل با استفاده از test set

۲. ثبت نام در سایت [kaggle.com](https://www.kaggle.com) و بارگذاری نتایج پیش‌بینی شده برای فایل `test.csv` و دریافت امتیاز (که همان دقت پیش‌بینی شما می‌باشد) از سایت [kaggle](https://www.kaggle.com). نحوه بارگذاری و فرمت فایل ارسالی برای سایت در این [لینک](#) توضیح داده شده است.

روش انتخابی شما برای محاسبه دقت باید در تمرین دوم و سوم یکسان باشد تا مقایسه بین درخت تصمیم و جنگل تصادفی به درستی صورت بگیرد.

**توجه ۲-** نیازی به مصورسازی جنگل تصادفی نیست.

## سوال ۹:

در این قسمت هدف آشنایی بیشتر و کلاس بندی با SVM (Support Vector Machine) می‌باشد. مجموعه داده مورد استفاده در اینجا همان مجموعه داده تایتانیك و هدف پیش‌بینی زنده ماندن یا نماندن مسافر خواهد بود. برای پیاده‌سازی این قسمت می‌توانید از کتابخانه‌های موجود استفاده کنید. مراحل زیر را انجام داده و نتایج خواسته شده را در گزارش خود بیاورید:

۱- پیش‌پردازش مناسب روی مجموعه داده انجام دهید. دقت کنید که نحوه ارزیابی مدل به یکی از دو روش ذکر شده در قسمت قبل خواهد بود (تقسیم داده‌های `train.csv` به دو قسمت `train set` و `test set` و یا استفاده از سایت [kaggle](https://www.kaggle.com))، لذا در صورت استفاده از روش اول داده‌های `train.csv` را در این مرحله به هم زده و تقسیم کنید.

۲- مدل خود را با SVM با هسته خطی (Linear kernel function) آموزش داده و دقت آن را محاسبه کنید.

۳- دقت به دست آمده در بخش قبل را با دقت درخت تصمیم (در تمرین قبل) و جنگل تصادفی مقایسه کنید. علت پایین‌تر (یا بالاتر) بودن دقت را از نظر خود بیان کنید.

۴- مدل خود را با SVM با یک هسته غیر خطی (به انتخاب خودتان) آموزش داده و دقت آن را به دست آورید.

۵- دقت پیش‌بینی در بخش ۲ و ۴ را با هم مقایسه کرده و علت تغییر آن را ذکر کنید.