

به نام خدا



دانشکده‌ی مهندسی کامپیوتر

سررسید تئوری: ۱۳ دی ماه چهارشنبه ۲۳:۵۹

سررسید عملی: ۱۵ دی ماه جمعه ۲۳:۵۹

پاییز ۱۴۰۲

## یادگیری ماشین

تمرین ۵: Adaboost, Decision Tree & Learning Principles

مدرس: مهدی جعفری سیاوشانی

- سررسید بخش تئوری این تمرین چهارشنبه ۱۳ دی ماه ساعت ۲۳:۵۹ است.
- سررسید بخش عملی این تمرین جمعه ۱۵ دی ماه ساعت ۲۳:۵۹ است.
- در صورت کشف تقلب، بار اول برای افراد درگیر تقلب، نمره‌ی همان سوال(های) خاص صفر در نظر گرفته می‌شوند.
- در صورت تکرار، نمره کل تمرین صفر در نظر گرفته می‌شود و در صورت تکرار، درس برای افراد حذف خواهد شد.
- تمامی پاسخ‌های خود را در یک فایل با فرمت (pdf) HW5-[SID]-[Fullname].zip روی کوئرا قرار دهید.

## پرسش‌ها

### ۱ قسمت تئوری

#### ۱.۱ پرسش اول (۱۰ نمره)

فرض کنید یک بانک اعتباری روشی خوبی برای اعتبار دادن به کاربران جدید ندارد. با رسیدن مشتری‌های  $x_1, x_2, \dots, x_N$  این بانک با روش اولیه خود یک سری از این مشتری‌های را تایید می‌کند و پس از گرفتن کارت اعتباری عملکرد آن‌ها را بررسی می‌کند. برای سادگی فرض کنید که به اولین  $N$  مشتری کارت اعتباری داده شده است. حال بانک برای بهبود الگوریتم خودش اطلاعاتی را که تا به حال جمع کرده به شما می‌دهد. این اطلاعات به صورت دوتایی‌های  $(x_1, y_1), \dots, (x_N, y_N)$  می‌باشد. شما پیش از حتی گرفتن اطلاعات، با یک سری از فرمول‌های ریاضی یک تابع به آن‌ها ارائه می‌دهید که به صورت بی‌نقص کار می‌کند.

۱. (۲,۵ نمره) اندازه مجموعه فرضیه یا  $M$  چقدر است؟

۲. (۲,۵ نمره) با این اندازه  $M$  باند هافدینگ چه چیزی را در مورد احتمال این که برای  $N = 10000$  میزان خطای کارکرد واقعی کمتر از دو درصد است، می‌گوید؟

۳. (۲,۵ نمره) شما جوابتان را به بانک می‌دهید و به آن‌ها اطمینان می‌دهید که کارکرد بهتر از خطای دو درصد خواهد بود و اطمینان شما از جواب قسمت قبل بدست می‌آید. بانک از جواب شما برای اعتبار دادن به کاربران جدید استفاده می‌کند. متأسفانه، بیش از نیمی از کارت‌های اعتبار آن‌ها به تأخیر می‌افتند. دلایل/دلیل ممکن برای این اتفاق چیست؟

۴. (۲,۵ نمره) آیا راهی وجود دارد که بانک بتواند با اطمینان احتمالاتی شما از تابعان استفاده کند؟ (راهنمایی: جواب مثبت است)

#### ۲.۱ پرسش دوم (۲۰ نمره)

۱. (۲ نمره) فرض کنید با یک مسئله دسته‌بندی دو دسته‌ای مواجه هستیم و تعداد  $m$  ویژگی داریم که هر ویژگی دو مقدار ممکن را می‌تواند اختیار کند. در این صورت نشان دهید چند درخت تصمیم متفاوت می‌توانیم برای این مسئله داشته باشیم. (توجه داشته باشید که هر نود دو زیر نود دارد که در واقع مجموعه مقادیر به دو دسته تقسیم می‌شوند و عضویت در هر کدام از زیرشاخه‌های درخت را تشکیل می‌دهند)

۲. (۳ نمره) آیا درخت تصمیمی که به طور حریصانه و با کمک بهینه‌کردن معیارهایی نظیر Information Gain ساخته می‌شود، همیشه بهترین درخت ممکن است؟ علت این امر را نیز توضیح دهید.

۳. درخت‌های تصمیم در برخورد با ویژگی‌های چند مقداره از دو استراتژی استفاده می‌کنند. در استراتژی اول که Multiway Split نام دارد، هر نود تصمیم‌گیری می‌تواند  $k$  خروجی داشته باشد. در حالی که در استراتژی دوم که نامش Binary Split است، هر نود تصمیم‌گیری صرفاً دو خروجی می‌تواند داشته باشد و در نتیجه اگر یک ویژگی داشته باشیم که  $k$  مقدار ممکن را بتواند اختیار کند، در هر نود تصمیم‌گیری در این استراتژی تنها می‌توانیم شرطی را چک کنیم که پاسخ به آن دو مقدار بله یا خیر باشد. جدول را در نظر بگیرید و با توجه به آن به سوالات زیر پاسخ دهید.

حزب سیاسی	محل زندگی	نژاد	جنسیت	ردیف
دموکرات	کالیفرنیا	سفیدپوست	مرد	۱
دموکرات	کالیفرنیا	سفیدپوست	مرد	۲
جمهوری خواه	تگزاس	سفیدپوست	مرد	۳
جمهوری خواه	تگزاس	سیاه پوست	مرد	۴
دموکرات	اوهاйо	سیاه پوست	مرد	۵
جمهوری خواه	کالیفرنیا	سفیدپوست	زن	۶
جمهوری خواه	تگزاس	سفیدپوست	زن	۷
دموکرات	اوهاйо	سیاه پوست	زن	۸
دموکرات	کالیفرنیا	سیاه پوست	زن	۹
جمهوری خواه	اوهاйо	سیاه پوست	زن	۱۰

(آ) (۴ نمره) داده‌های جدول را در نظر بگیرید. حزب سیاسی برچسب هدف و بقیه ستون‌ها نیز ویژگی هستند. با استفاده از معیار Gini Impurity و استراتژی Multiway Split درخت تصمیم با عمق حداکثر دو را برای این داده‌ها به دست آورید.

(ب) (۴ نمره) حال با استفاده از معیار Gini Impurity و استراتژی Binary Split درخت تصمیم با عمق حداکثر دو را بر این داده‌ها به دست آورید.

(ج) (۳ نمره) این دو استراتژی را با یکدیگر مقایسه کنید و مزایا و معایب‌شان نسبت به هم را تحلیل کنید.

۴. (۴ نمره) یکی از مشکلات درخت تصمیم این است که در عین حالی که به دقت بالایی دست پیدا می‌کند اما واریانس خطای آن بالاست. توضیح دهید که جنگل تصادفی چگونه با وجود حفظ دقت بالای درخت تصمیم، واریانس خطای آن را کاهش می‌دهد.

### ۳.۱ پرسش سوم (۲۰ نمره)

در این سوال نشان می‌دهید که شیوه انتخاب پارامتر  $\alpha_t$  در الگوریتم Adaboost معادل آن است که یک باند بالای نمایی برای تابع هزینه این الگوریتم را در هر دور به شکل حریصانه کمینه کنیم.

۱. (۴ نمره) فرض کنید  $h_t : \mathbb{R}^m \rightarrow \{-1, 1\}$  دسته‌بند ضعیفی است که در گام  $t$  ام به دست آورده‌ایم و  $\alpha_t$  وزن این دسته‌بند است. دیده‌ایم که دسته‌بند نهایی در الگوریتم Adaboost به شکل زیر است:

$$\hat{y} = \text{sign}(H_t(x)) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

که  $H_t$  دسته‌بند نهایی در پایان گام  $t$  است. فرض کنید  $\{(x_1, y_1), \dots, (x_N, y_N)\} \subset \mathbb{R}^m \times \{-1, 1\}$  داده‌های آموزش ما باشند. نشان دهید خطای دسته‌بند نهایی می‌تواند توسط یک تابع هزینه نمایی از بالا محدود شود:

$$E = \frac{1}{N} \sum_{i=1}^N \exp(-y_i H_t(x_i)) \geq \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{y}_i \neq y_i)$$

که  $\mathbb{I}$  تابع indicator است.

۲. (۳ نمره) فرض کنید  $D_t(i)$  توزیع وزن داده‌های آموزش در مرحله  $t$  است.  $D_{t+1}(i)$  را بر حسب  $y_i, x_i, \alpha_t, Z_t$  و دسته‌بند  $h_t$  بنویسید.  $T$  گام نهایی است و  $t \in \{1, \dots, T\}$ . به یاد داشته باشید که در  $Z_t$  عامل نرمال‌کننده توزیع  $D_{t+1}$  است:

$$Z_t = \sum_{i=1}^N D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

۳. (۳ نمره) نشان دهید که

$$E = \sum_{i=1}^N \frac{1}{N} \exp\left(\sum_{t=1}^T -\alpha_t y_i h_t(x_i)\right)$$

۴. (۴ نمره) نشان دهید که

$$E = \prod_{t=1}^N Z_t$$

۵. (۳ نمره) نشان دهید که  $Z_t$  می‌تواند به شکل زیر نوشته شود:

$$Z_t = (1 - \epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$$

که  $\epsilon_t$  خطای وزن‌دار دسته‌بند  $h_t$  است:

$$\epsilon_t = \sum_{i=1}^N D_t(i) \mathbb{I}(h_t(x_i) \neq y_i)$$

۶. (۳ نمره) همه گام‌های بالا را به آن دلیل طی کردیم که کمینه کردن خطای ۰ یا ۱ دسته‌بند روی داده‌های آموزش دشوار است. اما کمینه کردن حریصانه باند بالای  $E$  خطا، ممکن است. نشان دهید که انتخاب حریصانه  $\alpha_t$  در هر گام برای کمینه کردن  $Z_t$ ، به عبارت زیر منجر می‌شود:

$$\alpha_t = \frac{1}{\gamma} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$

## ۲ قسمت عملی

### ۱.۲ پرسش اول

در این سوال شما الگوریتم Adaboost را با دسته‌بندی‌های پایه decision stump پیاده‌سازی می‌کنید تا عمل دسته‌بندی روی مجموعه داده ساختگی adaboost-syndata انجام دهید. تولید داده مصنوعی (adaboost-syndata) به عهده خودتان است. مثلاً می‌توانید تابعی تعریف کنید که به آن تعداد نمونه مورد نیاز را بدهید و داده برچسب‌دار در فضای دو بعدی خروجی بدهد (به طور رندوم یا هر طریقی که برای مساله مناسب‌تر می‌دانید). یک decision stump یک جداکننده خطی است که موازی یکی از محورهاست و نمونه‌های یک سویش را مثبت و نمونه‌های سوی دیگر را منفی دسته‌بندی می‌کند. مثلاً فرض کنید ورودی‌ها دو بعدی هستند. در این صورت از decision stump های زیر استفاده می‌کنیم:

$$h(x) = \begin{cases} y & x_1 \geq k \\ -y & o.w. \end{cases}$$

یا

$$h(x) = \begin{cases} y & x_1 \leq k \\ -y & o.w. \end{cases}$$

که  $y \in \{1, -1\}$  و  $k$  یک عدد حقیقی دلخواه است.

طراحی کد شما بر عهده خودتان است اما باید یک تابع  $\text{Adaboost}(X_{\text{train}}, y_{\text{train}}, X_{\text{test}}, \text{num-iter})$  پیاده‌سازی کنید که یک بردار  $y_{\text{pred}}$  خروجی می‌دهد.  $\text{num-iter}$  تعداد دورهایی است که الگوریتم شما باید اجرا شود. دقت داشته باشید که داده‌ها دو بعدی هستند. بهتر است یک تابع داشته باشید که نمونه‌های یادگیری و وزن آن‌ها را ورودی می‌گیرد و خروجی بهترین decision stump را با توجه به وزن‌های فعلی و معیار decision stump و مکان آن، باز می‌گرداند. همچنین داشتن یک تابع برای به‌روزرسانی وزن‌ها و یک تابع برای پیش‌بینی برچسب‌ها با استفاده از decision stump های مراحل قبل توصیه می‌شود.

۱. (۷ نمره) در هر دور اجرای الگوریتم، خطای وزن‌دار  $\epsilon_t$  دسته‌بند ضعیف خود  $h_t$  را روی داده‌های آموزش محاسبه کنید. برای  $T = 20$  نمودار  $\epsilon_t$  را بر حسب  $t$  رسم کنید. بیشترین خطای یک دسته‌بند چقدر است؟ کجا رخ می‌دهد؟ نمودار را تفسیر کنید.

۲. (۶ نمره) در هر دور نرخ خطای دسته‌بند نهایی روی داده‌های آموزش  $H_t$  را محاسبه کنید. در  $T = 20$  نمودار این خطا را بر حسب  $t$  رسم کنید. نمودار را تفسیر کنید.

۳. (۶ نمره) در هر دور نرخ خطای دسته‌بند نهایی روی داده‌های تست  $H_t$  را محاسبه کنید. در  $T = 20$  نمودار این خطا را بر حسب  $t$  روی همان نمودار قسمت قبل رسم کنید. نمودار را تفسیر کنید.

۴. (۶ نمره) برای مقادیر  $T = [20, 50, 100, 200, 500, 1000, 2000, 4000]$  نرخ خطای آموزش و تست دسته‌بند نهایی را بیابید. نمودار این خطاها را رسم کنید. این نمودارها را تفسیر کنید.

## ۲.۲ پرسش دوم

داده‌های مورد استفاده در این تمرین مربوط به جمعیتی از قارچ‌ها است. از هر قارچ ۲۲ ویژگی استخراج شده است و هدف پیش‌بینی برچسب داده‌هاست. هدف طراحی دست‌بند مناسب برای این دادگان است. پیش از شروع دادگان را به سه بخش ۷۰ درصدی آموزش و ۲۰ درصدی اعتبار سنجی و ۱۰ درصدی سنجش تقسیم کنید. دقت داشته باشید که در این تمرین می‌توانید از کتابخانه `scikit-learn` استفاده نمایید.

۱. (۶ نمره) یکی از چالش‌های این دادگان نمونه‌هایی هستند که برخی از ویژگی‌های آن‌ها `miss` شده است. با این پدیده با استراتژی‌های متفاوتی می‌توان برخورد کرد. استراتژی مورد نظر خود در برخورد با این نمونه در این مسئله را در گزارش خود ارائه دهید.

۲. (۷ نمره) با استفاده از درخت تصمیم و با مقادیر حداکثر عمق  $\{4, 8, 16, 24, 32\}$  دسته‌بند مناسب طراحی کنید. خطای آموزش و تست را گزارش کنید. (معیار انتخاب ویژگی و تقسیم در هر گره را معیار `gini` در نظر بگیرید.)

۳. (۱۲ نمره) حال با استفاده از جنگل تصادفی یک دسته‌بند مناسب طراحی کنید. دسته‌بندی را با استفاده از ترکیب نتیجه ۷ عدد درخت انجام دهید. بهترین هایپرپارامتر برای هر یک از درخت‌های جنگل را از میان  $\max\_depth \in \{3, 5, 7\}$  و  $feature\_numbers \in \{3, 5, 7\}$  پیدا کنید. در نهایت دقت آموزش هر یک از تنظیمات جست‌وجو و دقت و صحت تست بهترین مدل انتخاب شده خود را گزارش کنید. (کلاس با برچسب `e` را به عنوان کلاس مثبت در نظر بگیرید)

موفق باشید