



دانشگاه صنعتی شریف

دانشکده‌ی مهندسی کامپیوتر

یادگیری ماشین

پاییز و زمستان ۱۴۰۴

استاد: علی شریفی زارچی

گردآورندگان: ارشیا یوسف‌نیا - بنیامین قبری - زهرا رحمانی - علی باوفا - کیهان هدایی

مهلت ارسال: ۳۰ آبان

یادگیری با ناظارت

تمرین اول

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.

- در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر مجاز داشته باشید و در مجموع حداکثر ۱۵ روز تأخیر مجاز خواهد داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت جداگانه محاسبه می‌شود و مجموع تأخیر هر دو باید بیشتر از ۱۵ روز شود. پس از اتمام زمان مجاز، دو روز اضافی برای آپلود غیرمجاز در نظر گرفته شده است که در این بازه، به ازای هر ساعت تأخیر، ۲ درصد از نمره تمرین کسر خواهد شد.
- همکاری و همفکری شما در انجام تمرین مانع ندارد اما پاسخ ارسالی هر کس حتماً باید توسط خود او نوشته شده باشد.
- در صورت همفکری و یا استفاده از هر منابع خارج درسی، نام همفکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفاً تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

سوالات نظری (۱۰۰ نمره)

۱. (۱۰ نمره) درستی و نادرستی گزاره‌های زیر را همراه یک استدلال کوتاه بیان کنید.

- (الف) اگر داده‌ها خطی قابل تفکیک نباشند، هیچ راه حلی برای SVM به صورت soft-margin نداریم.
- (ب) برای یک مجموعه داده که نویز زیادی دارد یعنی تعداد زیادی از دادگان برچسب نادرست دارند، استفاده از جنگل تصادفی به طور کلی بهتر از boosted decision trees است.
- (پ) در روش bagging در استفاده از درخت تصمیم‌گیری، معمولاً n درخت که واریانس زیاد و بایاس کمی دارند انتخاب می‌کنیم و نتیجه نهایی با یک اجماع از نتیجه آن‌ها به دست می‌آید.
- (ت) در تجزیه واریانس-بایاس، استفاده از یک مدل بسیار پیچیده بر روی داده محدود، معمولاً در ابتدا به بایاس کم و واریانس بالا منجر می‌شود.
- (ث) الگوریتم پرسپترون تضمین می‌کند که برای هر مجموعه داده‌ای که به صورت خطی قابل تفکیک باشند، پس از تعداد محدودی گام به یک راه حل همگرا می‌شود.

۲. (۱۵ نمره) روش تخمین چگالی در k nn برآورده از چگالی یک نقطه از فضای ویژگی‌ها را نسبت به چگالی k همسایه نزدیک‌تر می‌دهد. نشان دهید که مدل چگالی k nn را می‌توان یک توزیع نامناسب دید. این موضوع چه اهمیتی دارد و چه نتیجه‌ای می‌دهد؟

توزیع نامناسب (improper distribution) توزیعی است که انتگرال‌گیری روی آن روی کل بازه واگرا خواهد شد.

راهنمایی: برای محاسبه $(x_i)_i \in \mathbb{P}$, کره‌ای کوچک با مرکزیت خود آن نقطه در نظر بگیرید و شعاع را تا دربرگرفتن k نقطه افزایش دهید. افروزن بر این، کره را چنان کوچک فرض کنید که $(x_i)_i \in \mathbb{P}$ در آن به تقریب ثابت باشد.

۳. (۱۵ نمره) فرض کنید یکتابع رگرسیون $h(\mathbf{x})$ داریم که برای هر بردار ورودی \mathbf{x} برچسب $y = h(\mathbf{x})$ را به آن نسبت می‌دهد. از دادگان آموزش استفاده کردایم و m مدل پیش‌بینی‌کننده آموزش داده‌ایم که آنها را با $\hat{h}_1(\mathbf{x}), \hat{h}_2(\mathbf{x}), \dots, \hat{h}_m(\mathbf{x})$ نشان می‌دهیم. اکنون قرار دهید:

$$\hat{H}_m(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \hat{h}_i(\mathbf{x})$$

ثابت کنید که:

$$\mathbb{E}_{\mathbf{x}} \left[(\hat{H}_m(\mathbf{x}) - h(\mathbf{x}))^2 \right] \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{x}} \left[(\hat{h}_i(\mathbf{x}) - h(\mathbf{x}))^2 \right]$$

مفهوم این حکم چیست؟

۴. (۱۰ نمره) به سوالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(الف) یک مدل Random Forest روی تعداد زیادی داده برای یک مسئله دسته‌بندی آموزش داده‌ایم. میزان خطای مدل روی دادگان آموزشی پایین است ولی روی دادگان آزمون خطای زیادی گزارش شده است. دو دلیل برای این مشکل بیان کنید و برای هر کدام راه حلی بگویید.

(ب) فرض کنید دادگان $X = \{\mathbf{x}_i\}_{i=1}^n$ با میانگین و واریانس μ, σ^2 داده شده‌اند. بگیرید:

$$z_m = \frac{1}{m} \left(\sum_{i=1}^m a_i \right)$$

که در آن a_i ها تصادفی و متمایز از X گزینش شده‌اند. میانگین و واریانس z_m را بیابید. اکنون بیان کنید مجموعه دادگان $\{z_m\}$ چه خوبی و بدی‌هایی نسبت به X دارد و در چه شرایطی استفاده از آن برای آموزش مناسب است.

۵. (۱۰ نمره) به سوالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(الف) با توجه به تابع سیگموید به سوالات زیر پاسخ دهید.

- مشتق تابع سیگموید را بدست آورید (تمامی مراحل به طور کامل نوشته شود).
- تابع سیگموید و مشتق آن را رسم کنید. از نمودار بدست آمده کمک بگیرید و به طور شهودی توضیح دهید که مشکل تابع سیگموید در بروز رسانی وزن‌ها با کاهش گرادیان چیست.

(ب) چرا رگرسیون لجستیک که با الگوریتم کاهش گرادیان آموزش داده می‌شود، بعضی وقت‌ها که دادگان تقریباً به طور خطی قابل تفکیک هستند، کند همگرا می‌شود؟

(ج) در یک مسئله طبقه‌بندی چندکلاسه با K کلاس، به جای استفاده از چندین مدل دودویی، می‌توان از یک مدل چندکلاسه استفاده کرد. نشان دهید که اگر به همه بردارهای پارامتر $\theta_1, \theta_2, \dots, \theta_K$ یک بردار ثابت c اضافه کنیم، مقادیر احتمالات تغییری نمی‌کنند.

۶. (۱۵ نمره) به سوالاتی که در ادامه آمده‌اند با بیان توضیحات کامل پاسخ دهید.

(الف) برای یک مجموعه binary classification به صورت (\mathbf{x}_i, y_i) که $\{+1, -1\}$ مسئله-
margin^۴ را در نظر بگیرید که به صورت مسئله بهینه‌سازی

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\} \quad s.t. \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$

تعريف می‌شود.

- توضیح دهید چرا Geometric Margin مربوط به این classifier برابر با

$$\gamma = \min_i \frac{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

است.

- چرا با کاهش دادن $\frac{1}{2} \|\mathbf{w}\|^2$ با توجه به شروط داده شده، این margin افزایش می‌یابد؟

(ب) فرض کنید مدل Linear Regression را داریم:

$$\mathbf{y} = \mathbf{X}\beta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

- فرم بسته‌ای برای تخمین‌گر least squares برای پارامتر β به دست آورید.

- نشان دهید که این تخمین‌گر unbiased است و ماتریس covariance آن را نیز به دست آورید.

۷. (۱۰ نمره) همان‌طور که می‌دانید اعتبارسنجی مقاطعه^۳ با حذف یک نمونه یا LOOCV را می‌توان بدین گونه نشان داد:

$$e(S_n) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{-i}(\mathbf{x}_i))^2$$

که در آن $S_n^{-i} = S_n \setminus \{(x_i, y_i)\}$ (مجموعه داده‌های آموزش)، $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ (مجموعه داده‌ها پس از حذف نمونه iام)، و \hat{f}_{-i} مدل آموزش‌دیده روی S_n^{-i} می‌باشد.

(الف) با فرض اینکه فقط نمونه اول حذف شده است، خطاب به شکل

$$e_1(S_n) = (y_1 - \hat{f}_{-1}(\mathbf{x}_1))^2$$

است. نشان دهید:

$$\mathbb{E}[e_1(S_n)] = \mathbb{E}[(y - \hat{f}_{S_{n-1}}(\mathbf{x}))^2]$$

که در آن (y, \mathbf{x}) یک نمونه تصادفی از توزیع داده‌های $\hat{f}_{S_{n-1}}$ است.

(ب) با استفاده از نتیجه بالا، نشان دهید که:

$$\mathbb{E}[e(S_n)] = \mathbb{E}[(y - \hat{f}_{S_{n-1}}(\mathbf{x}))^2]$$

۸. (۱۵ نمره) جدول دادگان زیر را برای یک مسئله دسته‌بندی غذا در نظر بگیرید:

دما	اندازه	مزه	خواص‌یابنده
گرم	کوچک	شور	خیر
سرد	بزرگ	شیرین	خیر
سرد	بزرگ	شیرین	خیر
سرد	کوچک	ترش	بله
گرم	کوچک	ترش	بله
گرم	بزرگ	شور	خیر
گرم	بزرگ	ترش	بله
سرد	کوچک	شیرین	بله
سرد	کوچک	شیرین	بله
گرم	بزرگ	شور	خیر

(الف) آنتروپی^۴ اولیه ستون مربوط به خواص‌یابنده بودن غذا چند است؟

[https://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics))^۳

<https://en.wikipedia.org/wiki/Entropy>^۴

(ب) فرض کنید ویژگی مزه را ریشه درخت گرفته‌ایم، اکنون gain information^۵ مربوط به این ویژگی را مشخص کنید.

(پ) درخت تصمیم کامل را برای اجرای پیش‌بینی رسم کنید و همه محاسبات و مراحل را بیاورید.

(ت) اگر یک غذا شیرین، کوچک، و سرد باشد، پیش‌بینی ما با درخت قسمت قبل چه خواهد بود؟

(ث) یک نمونه bootstrap^۶ از دادگان جدول گرفته و درخت تصمیم‌گیری را دوباره بسازید.