



هوش مصنوعی

بهار ۱۴۰۲

استاد: محمدحسین رهبان

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

گردآورندگان: محمدرضا دویران، امیررضا میرزایی و محمد جواد هزاره

آشنایی با یادگیری ماشین، رگرسیون، درخت تصمیم‌گیری مهلت ارسال: ۱۹ خرداد

تمرین پنجم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمرین تا سقف ۷ روز و در مجموع ۱۵ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر روز تأخیر غیر مجاز ۱۲ درصد از نمره تمرین به صورت ساعتی کسر خواهد شد.
- هم‌کاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ‌های ارسال هر کس حتما 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.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

سوالات نظری (۱۲۰ نمره)

۱. (۲۰ نمره) با توجه به مفاهیم فراگرفته شده در درس به سوالات زیر پاسخ دهید.
(الف) برای مدل‌هایی که اریبی زیادی دارند دست کم دو راهکار ارائه دهید که مقدار این اریبی کاهش یابد.
(ب) فرض کنید تعدادی از ویژگی‌های مدل دوبه‌دو با یکدیگر هم‌بسته باشند. این اتفاق را از دیدگاه bias-variance بررسی کنید. اگر ویژگی‌های هم‌بسته را حذف کنیم بیان کنید که چگونه bias و variance تغییر می‌کنند.
(پ) کدام یک از گرازه‌های زیر درست می‌باشند؟ چرا؟
 - اگر مقدار بایاس زیاد باشد، افزایش تعداد داده‌های آموزش می‌تواند باعث کاهش بایاس شود.
 - افزایش پیچیدگی مدل در رگرسیون همواره باعث کاهش خطای آموزش و افزایش خطای تست می‌شود.
۲. (۳۰ نمره)
(آ) فرض کنید برای داده‌های جدول ۱ یک درخت تصمیم آموزش می‌دهیم تا X را به وسیله A, B, C پیش‌بینی کنیم. درصد خطای مدل پس از آموزش بر روی داده‌های آموزش چقدر خواهد بود؟
(ب) فرض کنید روی مجموعه‌ی داده‌ی دلخواهی، درخت تصمیمی برای دسته‌بندی بین k کلاس، آموزش می‌دهیم. حداکثر خطایی که ممکن است این مدل روی داده‌های آموزش داشته باشد چقدر خواهد بود؟ (پاسخ را به صورت کسری بنویسید)

C	B	A	X
۰	۰	۰	۰
۱	۰	۰	۰
۱	۰	۰	۰
۰	۱	۰	۰
۱	۱	۰	۰
۱	۱	۰	۱
۱	۱	۰	۱
۰	۰	۱	۰
۱	۰	۱	۱
۰	۱	۱	۱
۰	۱	۱	۱
۱	۱	۱	۰
۱	۱	۱	۱

جدول ۱: داده‌های مدل درخت تصمیم

۳. (۴۰ نمره) در رابطه با الگوریتم Logistic regression به سوالات زیر پاسخ دهید.

الف) این الگوریتم را برای حالت K کلاسه تغییر دهید و احتمالات آن را بنویسید.

ب) همانطور که در قسمت الف به دست آوردید، در Logistic regression برای K کلاس، احتمال پسین به روش زیر محاسبه می‌شود:

$$P(Y = k|X = x) = \frac{e^{w_k^T x}}{1 + \sum_{i=1}^{K-1} e^{w_i^T x}}, k = 1, 2, 3, 4, \dots, K-1$$

$$P(Y = K|X = x) = \frac{1}{1 + \sum_{i=1}^{K-1} e^{w_i^T x}}$$

برای راحتی فرض کردیم که $w_k = 0$ می‌باشد. کدام یک از پارامترها باید تخمین زده شوند؟

پ) حال log-likelihood زیر را برای n نمونه‌ی زیر ساده کنید:

$$\text{Samples} : (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

$$L(w_1, \dots, w_{K-1}) = \sum_{i=1}^n \ln P(Y = y_i | X = x_i)$$

ت) گرادیان L را نسبت به هریک از w_k ها بیاید و آن را ساده کنید.

ث) تابع هدف زیر را در نظر بگیرید. گرادیان f را با توجه به هریک از w_k ها بیاید.

$$f(w_1, \dots, w_{K-1}) = L(w_1, \dots, w_{K-1}) - \frac{\lambda}{2} \sum_{j=1}^{K-1} \|w_j\|_2^2$$

۴. (۳۰ نمره) فرض کنید n داده آموزش با m ویژگی داریم که که ماتریس این داده‌ها را $X_{n \times m}$ در نظر می‌گیریم.

بردار مقدار هدف نیز برابر $y = [y^{(1)}, \dots, y^{(n)}]$ می‌باشد. در ادامه منظور از x_j , j امین ستون ماتریس X است. حال با توجه به توضیحات داده شده به سوالات زیر پاسخ دهید.

الف) ابتدا ثابت کنید اگر رگرسیون را فقط بر روی یکی از m ویژگی موجود آموزش دهیم آنگاه خواهیم داشت:

$$w_j = \frac{x_j^T y}{x_j^T x_j}$$

ب) فرض کنید ستون‌های ماتریس X متعامد باشد. ثابت کنید که پارامترهای بهینه از آموزش رگرسیون بر روی همه ویژگی‌ها با پارامترهای بهینه حاصل از آموزش روی هر ویژگی به طور مستقل یکسان است.

پ) فرض کنید می‌خواهیم یک رگرسیون بر روی بایاس و یکی از ویژگی‌های نمونه داده‌ها آموزش دهیم. $(w = [w_j, w_0])$ با توجه به اطلاعات داده شده عبارات زیر را اثبات کنید:

$$w_j = \frac{\text{cov}[x_j, y]}{\text{var}[x_j]}$$

$$w_0 = E[y] - w_j E[x_j]$$

سوالات عملی (۲۰۰ نمره)

۱. (۱۰۰ نمره) برای حل سوال عملی اول به نوت‌بوک `decision_tree.ipynb` مراجعه کنید.
۲. (۱۰۰ نمره) برای حل سوال عملی دوم به نوت‌بوک `logistic_regression.ipynb` مراجعه کنید.