



## مفاهیم پیشرفته در یادگیری ماشین

نیم سال دوم ۱۴۰۱-۱۴۰۰

مدرس: دکتر مهدیه سلیمانی

زمان تحویل: ۲۵ اسفندماه

عنوان تمرین

تمرین سری اول

لطفا نکات زیر را رعایت کنید:

- سوالات خود را از طریق پست مربوط به تمرین در Quera مطرح کنید.
- در هر کدام از سوالات، اگر از منابع خاصی استفاده کرده‌اید باید آن را ذکر کنید.
- پاسخ ارسالی واضح و خوانا باشد.
- تمام پاسخ‌های خود را در یک فایل با فرمت [Fullname].zip\_[SID]\_[HW#] روی کوثر قرار دهید.
- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. علاوه بر آن، در هر تمرین می‌توانید تا سقف هفت روز از تأخیر مجاز باقیمانده خود استفاده کنید.
- برای کسب نمره کامل در این تمرین کفایت ۱۶۰ نمره را دریافت نمایید، ما بقی نمرات امتیازی می‌باشند (۴۰ نمره امتیازی).

## سوال ۱: بهینه‌سازی در یادگیری چند وظیفه‌ای (۱۵ نمره)

مشکلی که متایادگیری مبتنی بر پروتوتایپ دارند این است که دسته‌بند ساده‌ای دارند و تعمیم‌پذیری دسته‌بند در فاز Meta-test کم می‌باشد. به همین دلیل در مقاله Bertinetto ۲۰۱۸ به بررسی دو دسته‌بند رایج در یادگیری ماشین (Logistic Regression, Ridge Regression) و نحوه استفاده موثر آن در مسئله متایادگیری پرداخته است. با مطالعه مقاله و راهنمایی‌های داده شده در زیر، به سوالات پاسخ دهید

(آ) نحوه استفاده از Ridge Regression برای دسته‌بندی را به طور کامل شرح دهید.

(ب) دو ماتریس  $X \in n \times d$  و  $Y \in n \times l$  را در نظر بگیرید که  $n$  تعداد داده‌ها و  $d$  اندازه بردار بازنمایی هر داده و  $l$  تعداد برچسب‌های موجود در دسته می‌باشد. ثابت کنید که دو رابطه زیر در صورتی که  $\lambda > 0$  باشد، با هم برابر می‌باشند:

$$(X^T X + \lambda I)^{-1} X^T Y = X^T (X X^T + \lambda I)^{-1} Y$$

و بیان کنید که در ادبیات متایادگیری استفاده از کدام یک از دو رابطه بالا در حل مسئله بهتر می‌باشد و چرا؟

(ج) در روش Logistic Regression برای بدست آوردن دسته‌بند برخلاف روش قبلی نیازمند بهینه‌سازی می‌باشیم. در مقاله از Newton's Method برای بهینه‌سازی استفاده شده است. دلیل این امر را بیان کنید و همچنین در مورد خود Newton's Method تحقیق کنید و رابطه به روزرسانی و نحوه بدست آوردن این رابطه را بنویسید.

(د) با اعمال Newton's Method روی تابع هزینه این دسته‌بند، رابطه شماره ۷ مقاله که به‌روزرسانی وزن‌های مدل می‌باشد را بدست آورید.

پاسخ:

(آ) طبق توضیحات مقاله، تصویر ورودی توسط یک Feature Extractor به یک بردار نگاشت می‌شود و سپس با استفاده از ماژول Regressor Ridge پیش‌بینی مدل از روی فضای بازنمایی بدست می‌آید. مشکلی که ممکن است با آن مواجه شویم این است که خروجی مدل Regressor مناسب تابع هزینه دسته‌بندی نباشد. به همین یک تبدیل affine روی خروجی regressor اعمال می‌شود که پارامترهای آن می‌تواند به عنوان متاپارامتر در حلقه بیرونی الگوریتم متایادگیری بدست آید.

(ب) باتوجه به اینکه ماتریس  $Y$  در هر دو طرف مساوی از سمت راست در معادله ضرب شده است، پس تنها کافی است که اثبات کنیم:

$$(X^T X + \lambda I)^{-1} X^T = X^T (X X^T + \lambda I)^{-1}$$

برای شروع اثبات در رابطه زیر را در نظر بگیرید:

$$\lambda X^T = \lambda X^T$$

ماتریس همانی را برای یک سمت از رابطه مساوی بالا از سمت چپ ماتریس  $X^T$  و یک بار از سمت راست ماتریس، ضرب می‌کنیم:

$$\lambda I_d X^T = \lambda X^T I_n$$

حال عبارت  $X^T X X^T$  را به دو سمت مساوی اضافه می‌کنیم.

$$X^T X X^T + \lambda I_d X^T = X^T X X^T + \lambda X^T I_n$$

با فاکتورگیری داریم:

$$(X^T X + \lambda I_d) X^T = X^T (X X^T + \lambda I_n)$$

حال اگر در دو طرف تساوی، عبارت  $(X^T X + \lambda I_d)^{-1}$  را از چپ و عبارت  $(X X^T + \lambda I_n)^{-1}$  را از سمت راست ضرب کنیم داریم:

$$(X^T X + \lambda I_d)^{-1} (X^T X + \lambda I_d) X^T (X X^T + \lambda I_n)^{-1} = (X^T X + \lambda I_d)^{-1} X^T (X X^T + \lambda I_n) (X X^T + \lambda I_n)^{-1}$$

با ساده‌سازی داریم:

$$X^T (X X^T + \lambda I_n)^{-1} = (X^T X + \lambda I_d)^{-1} X^T$$

که مطلوب سوال می‌باشد.

این نکته لازم به ذکر می‌باشد که ماتریس‌هایی که معکوس آن را در اثبات بالا استفاده کردیم هردو ماتریس‌های مثبت نیمه‌معین می‌باشند و بنابراین حتماً معکوس‌پذیر می‌باشند.

در ادبیات متیادگیری استفاده از رابطه  $X^T (X X^T + \lambda I_n)^{-1}$  بهینه‌تر می‌باشد. چون در این رابطه نیاز به محاسبه معکوس یک ماتریس با ابعاد  $n \times n$  داریم ولی در حالت دیگر نیاز به محاسبه معکوس ماتریسی به ابعاد  $d \times d$  می‌باشد که مقدار  $d$  که بیانگر ابعاد بازنمایی در شبکه عصبی می‌باشد که به مراتب خیلی بزرگتر از تعداد نمونه‌ها در مسئله‌های متیادگیری می‌باشد که تعداد نمونه‌های آموزش خیلی کم می‌باشد.

(ج) به دلیل محدودیت تعداد به‌روزرسانی پارامترها در حلقه درونی الگوریتم، از Newton's Method برای همگرایی سریع‌تر به نسبت گرادینت گیری ساده استفاده شده است. در Newton's Method علاوه بر گرادینت، اطلاعات مشتق دوم نیز استفاده می‌شود. در این روش با نوشتن بسط Taylor حول یک نقطه داریم:

$$f(x+t) = f(x) + f'(x)t + \frac{f''(x)t^2}{2}$$

در رابطه بالا برای کمینه کردن مقدار تابع، نسبت به پارامتر  $t$  که جهت به‌روزرسانی و حرکت می‌باشد مشتق می‌گیریم که جهت بهینه به‌روزرسانی را پیدا کنیم:

$$\frac{dy}{dx} \left( f(x) + f'(x)t + \frac{f''(x)t^2}{2} \right) = f'(x) + f''(x)t = 0$$

$$t^* = -\frac{f'(x)}{f''(x)}$$

(د) هدف سوال بیشینه‌کردن بیشینه درست‌نمایی می‌باشد.

$$P(Y|\omega, X) = \prod_{i=1}^N (\sigma(\omega^T x_i))^{y_i} (1 - \sigma(\omega^T x_i))^{1-y_i}$$

$$L = -\prod_{i=1}^N (\sigma(\omega^T x_i))^{y_i} (1 - \sigma(\omega^T x_i))^{1-y_i} = -\sum_{i=1}^N [y_i \log(\sigma(\omega^T x_i)) + (1 - y_i) \log(1 - \sigma(\omega^T x_i))]$$

حال با گرادینت گیری مقادیر گرادینت اول و دوم را حساب می‌کنیم:

$$\nabla L = \sum_{i=1}^N (\sigma(\omega^T x_i) - y_i) x_i^T$$

$$\nabla \nabla L = \sum_{i=1}^N \sigma(\omega^T x_i) (1 - \sigma(\omega^T x_i)) x_i^T x_i$$

ابتدا چند ماتریس تعریف می‌کنیم:

$$A_t = \text{diag}(q_t)$$

$$q_t^{(i)} = \sigma(\omega_t^T x^{(i)}) (1 - \sigma(\omega_t^T x^{(i)}))$$

$$B_t^{(i)} = \sigma(\omega_t^T x^{(i)}) - y^{(i)}$$

که  $A \in n \times n$  یک ماتریس قطری و  $B \in n$  می‌باشد. حال اگر سیگماهای بالا را به ضرب ماتریسی تبدیل کنیم، طبق رابطه به‌روزرسانی داریم:

$$\begin{aligned}\omega_{t+1} &= \omega_t - H^{-1} \nabla L \\ &= \omega_t - (X^T A_t X + \lambda I)^{-1} (X^T B_t + \lambda \omega_t) \\ &= (X^T A_t X + \lambda I)^{-1} ((X^T A_t X + \lambda I) \omega_t - X^T B_t - \lambda \omega_t) \\ &= (X^T A_t X + \lambda I)^{-1} (X^T A_t X \omega_t - X^T B_t)\end{aligned}$$

### سوال ۲: (نظری) تنظیم توزیع برای یادگیری چندنمونه‌ای (۴۰ نمره)

یکی از ریسک‌های احتمالی در یادگیری چندنمونه‌ای احتمال بیش برآزش بر روی دادگان کم‌تعداد آموزشی است. در این مقاله روشی پیشنهاد شده است تا به کمک استخراج مشخصات آماری کلاس‌های حاضر در متآموزش بتوان توزیع دادگان کلاس‌های حاضر در متانتست را تنظیم کرد. این مقاله را به دقت خوانده و به سوالات زیر به طور کامل پاسخ دهید:

(آ) از آنجایی که ممکن است توزیع دادگان هر کلاس حاضر در متآموزش گاوسی نباشد و دارای مقداری کشیدگی باشد؛ در نظر گرفتن این توزیع‌ها به عنوان توزیع گاوسی و استخراج میانگین و کواریانس از آن‌ها می‌تواند اشتباه باشد. توضیح دهید این مقاله چه روشی را برای حل مشکل کشیدگی توزیع دادگان متآموزش اتخاذ کرده است و چگونه این روش موجب حل مشکل کشیدگی توزیع می‌شود؟

(ب) پس از استخراج میانگین و کواریانس کلاس‌های حاضر در متآموزش، مدل ارائه شده اقدام به تنظیم توزیع دادگان حاضر در متانتست می‌کند. به صورت کامل و با نوشتن روابط ریاضی مربوطه بیان کنید این تنظیم توزیع به چه صورت انجام می‌پذیرد و وجود کلاس‌های مشابه در متانتست به کلاس منظور در متانتست چه کمکی به تنظیم توزیع می‌کند؟

(ج) در تنظیماتی که در هنگام متانتست از هر کلاس بیش از یک نمونه آموزش داشته باشیم این مدل به جای میانگین‌گیری از نمونه‌ها، برای هر کدام از  $k$  نمونه آموزش اقدام به تنظیم توزیع جداگانه می‌کند. توضیح دهید توزیع تنظیم جداگانه چه مزیتی نسبت به میانگین‌گیری نمونه‌ها و سپس یک توزیع تنظیم دارد؟

### سوال ۳: (عملی) یادگیری چندنمونه‌ای از طریق یادگیری متر (۴۰ نمره)

در این سوال قصد داریم تا مدل یادگیرنده شبکه **Prototypical** را مورد پیاده‌سازی و بررسی قرار دهیم. به این منظور هر دو زیر مجموعه‌های آموزش و تست دادگان CIFAR100 را دریافت کرده و سپس آن‌ها را به یکدیگر الحاق نمایید. سپس این دادگان را به سه زیر مجموعه متآموزش، متااعتبارسنجی و متانتست تقسیم کنید. به این صورت که دادگان آموزش شامل دادگان ۷۰ کلاس، دادگان اعتبارسنجی شامل ۲۰ کلاس و دادگان کلاس تست شامل ۱۰ دیگر باشند. در گام بعدی بایستی یک **Sampler** پیاده‌سازی کنید که با گرفتن پارامترهای **Shot** و **Way** بتواند دادگان اتکا و پرسمان برای یک وظیفه را تولید کند (در واقع این ماژول با هر فراخوانی دو مجموعه داده به اندازه  $Way * Shot$  برای اتکا و پرسمان خروجی می‌دهد). این ماژول در واقع هر بار یک اپیزود را تولید می‌کند. در طی آزمایش‌های زیر از ماژول **ProtoNetBack** که در فایل‌های پیوستی قرار داده شده است به عنوان شبکه تکیه استخراج ویژگی استفاده کنید و در جلوی آن دو لایه متصل با اندازه دلخواه قرار دهید.

(آ) دسته‌بند را با تنظیمات 8-shot, 10-way آموزش دهید و سپس دقت مدل را بر روی دادگان متانتست گزارش دهید. انتظار می‌رود دقت در این بخش بیشتر از ۵۰ درصد باشد.

(ب) به ازای هر یک از تنظیمات  $shot \in \{1, 2, 4, 8, 16\}$  و 10-way آزمایش بالا را تکرار کرده و نمودار دقت متانتست بر حسب shot را رسم نمایید.

(ج) حال به ازای هر یک از تنظیمات  $way \in \{2, 4, 8, 16, 32\}$  و 5-shot آزمایش را تکرار کرده و نمودار دقت متانتست بر حسب way را رسم نمایید. (دقت کنید که در هنگام متانتست از تنظیمات 5-shot, 10-way استفاده نمایید)

(د) حال در هنگام متآموزش با تنظیمات 10-shot, 10-way دسته‌بند را آموزش دهید. در هنگام متانتست اما متانتست را به ازای هر یک از تنظیمات  $shot \in \{1, 5, 10, 15, 20\}$  انجام دهید و نمودار دقت آن را بر حسب shot رسم نمایید.

### سوال ۴: (عملی) متایادگیری براساس بهینه‌سازی (۴۰ نمره)

در این سوال قصد داریم تا مدل معروف دسته متایادگیری براساس بهینه‌سازی، **MAML**، را پیاده‌سازی نماییم. مقاله مرتبط با این کار، این مقاله می‌باشد. در Notebook داده شده تمام پارامترهای مسئله و مراحل حل به صورت گام به گام تشریح شده است. سوال از دو بخش اصلی تشکیل

شده است که در بخش اول به دلیل کاهش هزینه آموزش بخش عمده شبکه به صورت pretrained شده در اختیار شما قرار داده شده است و شما تنها روی بخش مشخص شده شبکه فرایند متیادگیری را انجام خواهید داد. در بخش اول قرار است تاثیر تعداد گام‌های به‌روزرسانی مدل در حلقه داخلی الگوریتم، مورد بررسی قرار گیرد. از شما خواسته شده است که به ازای مقادیر ۱ تا ۳ این مورد را انجام دهید و نتیجه هر حالت را مقایسه و گزارش کنید. در بخش دوم نیز از شما خواسته شده است که حال با یک گام به‌روزرسانی حلقه داخلی، کل ساختار مدل (مدل متیادگیری بخش اول + ساختار مدل pretrained داده شده) را به صورت متاپارامتر در نظر بگیرید و متیادگیری را روی آن انجام دهید. در نهایت نتایج بدست آمده از هر دو بخش را تحلیل و گزارش نمایید.