



# مفاهیم پیشرفته در یادگیری ماشین

نیم سال دوم ۱۴۰۰-۴۰۱  
مدرس: دکتر مهدیه سلیمانی

زمان تحویل: ۱۹ فروردین

روش های متا-یادگیری مبتنی بر بهینه سازی و یادگیری متریک

تمرین سری دوم

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 کرده اید باید آن را ذکر کنید.
- پاسخ ارسالی واضح و خوانا باشد.
- تمام پاسخ های خود را در یک فایل با فرمت [Fullname]\_[SID]\_HW# روی کوثر قرار دهید.
- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. علاوه بر آن، در هر تمرین می توانید تا سقف هفت روز از تأخیر مجاز باقیمانده خود استفاده کنید.
- برای کسب نمره کامل در این تمرین کفایت تا ۱۶۰ نمره از ۲۰۰ نمره را کسب نمایید (۴۰ نمره امتیازی).

## سوال ۱: (نظری) متا-یادگیری مبتنی بر بهینه سازی دو سطحی (۲۵ نمره)

همانطور که در جلسات درس مشاهده کرده اید، یکی از روش های متا-یادگیری<sup>۱</sup>، خانواده بهینه سازی دو سطحی<sup>۲</sup> بوده که مهمترین کار در این زمینه روش MAML<sup>۳</sup> می باشد. در این خانواده از روش ها، متا-پارامترها<sup>۴</sup> (پارامترهای آهسته) همبند با پارامترهای مختص وظیفه<sup>۵</sup> (پارامترهای سریع) بوده و به عنوان یک نقطه شروع برای کل وزن های شبکه عمل می کنند. به طور دقیق تر، اگر توزیع وظایف<sup>۶</sup> را به صورت  $p(\mathcal{T})$  در نظر بگیریم، می توان رابطه زیر را برای یادگیری متا-پارامترها ارائه داد:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\mathcal{T}=(S,Q) \sim p(\mathcal{T})} [\mathcal{L}(\phi, Q)] \quad (1)$$

$$\text{Where } \phi = \text{Alg}(\theta, S) \quad (1b)$$

که در این رابطه  $S$  مجموعه داده های پشتیبان<sup>۷</sup> و  $Q$  مجموعه داده های پرسمان<sup>۸</sup> مربوط به هر وظیفه<sup>۹</sup> را نشان می دهد. در این رابطه، محاسبه پارامترهای سریع  $\phi$  توسط روش Alg انجام می شود که در مقالات مربوطه به طرق مختلفی انتخاب شده و بهینه سازی داخلی<sup>۱۰</sup> نامیده می شود. همچنین در رابطه فوق، بهینه سازی خارجی<sup>۱۱</sup> که روی پارامترهای  $\theta$  صورت می پذیرد، به شکل  $\arg \min$  نمایش داده شده است. برای انجام بهینه سازی خارجی، لازم است تا از تابع معرفی شده نسبت به  $\theta$  گرادینان را به صورت زیر محاسبه کرده و  $\theta$  را از طریق آن به روزرسانی نماییم:

$$\nabla_{\theta} \mathbb{E}[\mathcal{L}(\phi, Q)] = \mathbb{E}[\nabla_{\theta} \mathcal{L}(\text{Alg}(\theta, S), Q)] \quad (2)$$

برای محاسبه عبارت فوق، از قاعده زنجیره ای مشتق استفاده می کنیم:

$$\nabla_{\theta} \mathcal{L}(\text{Alg}(\theta, S), Q) = \nabla_{\phi} \mathcal{L}(\phi, Q)|_{\phi=\text{Alg}(\theta, S)} \times \frac{d}{d\theta} \text{Alg}(\theta, S) \quad (3)$$

همانطور که مشاهده می کنید، رابطه فوق از دو جمله تشکیل شده است؛ محاسبه جمله اول نسبتاً راحت است. چرا که کفایت تا از  $\mathcal{L}$  مشتق گرفته و مقدار  $\text{Alg}(\theta, S)$  را در آن جایگذاری کنیم و در این صورت گرادینان از  $\text{Alg}(\theta, S)$  عبور نمی کند. این در حالیست که برای محاسبه جمله دوم،

<sup>1</sup>Meta-Learning

<sup>2</sup>Bi-Level Optimization

<sup>3</sup>Model Agnostic Meta-Learning

<sup>4</sup>Meta-Parameters

<sup>5</sup>Task-Specific Parameters

<sup>6</sup>Tasks

<sup>7</sup>Support

<sup>8</sup>Query

<sup>9</sup>Task

<sup>10</sup>Inner-Level Optimization

<sup>11</sup>Outer-Level Optimization

لازم است تا عملیات مشتق‌گیری را از داخل الگوریتم  $\text{Alg}(\theta, S)$  عبور دهیم. این مسئله می‌تواند مشکل‌زا باشد چرا که ممکن است  $\text{Alg}(\theta, S)$  اصلاً قابلیت عبور گرادیان را نداشته باشد یا اگر دارد ممکن است منجر به محاسبات پرهزینه شود. به عنوان مثال، در روش MAML داریم:

$$\text{Alg}(\theta, S) = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, Q) \Rightarrow \frac{d}{d\theta} \text{Alg}(\theta, S) = I - \alpha \nabla_{\theta}^2 \mathcal{L}(\theta, S) \quad (۴)$$

که در رابطه فوق، محاسبه مشتق‌های مرتبه دوم (Hessian) می‌تواند بسیار سخت یا پرهزینه باشد چرا که لازم است تا کل گراف محاسباتی در طول مسیر محاسبه  $\text{Alg}(\theta, S)$  نگه داشته شود تا بتوان گرادیان را از روی آن عبور داده و به مراحل قبلی رساند. در این سوال قصد داریم تا تکنیکی معرفی کنیم که این مشکلات مرتفع شوند. برای این منظور، فرض کنید  $\text{Alg}(\theta, S)$  به صورت زیر پیشنهاد شده است:

$$\phi = \text{Alg}(\theta, S) = \arg \min_{\phi'} \mathcal{L}(\phi', S) + \frac{\lambda}{2} \|\phi' - \theta\|^2 \quad (۵)$$

که  $\lambda$  یک هایپرپارامتر است و هر چه مقدار آن بیشتر باشد، باعث می‌شود تا جواب بهینه‌سازی درونی، به نقطه شروع خود یعنی  $\theta$  نزدیک‌تر بماند. بهینه‌سازی ۵ را می‌توان با انجام چندین گام بهینه‌سازی تکرار شونده<sup>۱۲</sup> حل کرد اما مشکل آن جاست که امکان عبور گرادیان از چنین محاسباتی وجود ندارد. با در نظر گرفتن این مسئله، به پرسش‌های زیر پاسخ دهید:

(آ) فرض کنید ما قادر هستیم تا بهینه‌سازی ۵ را به صورت کامل حل کنیم و  $\phi$  را به عنوان جواب بهینه دقیق آن به دست آورده‌ایم. از این مسئله استفاده کنید و  $\frac{d\phi}{d\theta}$  را محاسبه کنید (راهنمایی: در نقطه بهینه دقیق، مشتق  $\mathcal{L}(\phi', S) + \frac{\lambda}{2} \|\phi' - \theta\|^2$  نسبت به  $\phi'$  صفر خواهد بود). (۱۰ نمره)

(ب) اگر مراحل پرسش قبل را به درستی طی کرده باشید، در جواب خود به یک عبارت حاوی مشتق مرتبه دوم  $\nabla^2 \mathcal{L}$  (یا همان ماتریس Hessian) می‌رسید. محاسبه این عبارت چه تفاوتی با مشتق مرتبه دوم موجود در رابطه ۴ دارد؟ استفاده از این تابع  $\text{Alg}$  پیشنهادی چه مزیتی نسبت به MAML دارد؟ (۱۰ نمره)

(ج) به عنوان جمع‌بندی، الگوریتمی که متا-یادگیر در هر اپیزود طی می‌کند را به صورت گام گام شرح دهید. (۵ نمره)

پاسخ

(آ)

$$\frac{d}{d\phi'} \mathcal{L}(\phi', S) + \frac{\lambda}{2} \|\phi' - \theta\|^2 = \nabla_{\phi'} \mathcal{L}(\phi', S)|_{\phi'=\phi} - \lambda(\phi' - \theta)|_{\phi'=\phi} = 0 \quad (۶آ)$$

$$\Rightarrow \phi = \theta - \frac{1}{\lambda} \nabla_{\phi} \mathcal{L}(\phi, S) \quad (۶ب)$$

$$\Rightarrow \frac{d}{d\theta} \phi = I - \frac{1}{\lambda} \nabla_{\phi}^2 \mathcal{L}(\phi, S) \times \frac{d\phi}{d\theta} \quad (۶ج)$$

$$\Rightarrow \frac{d\phi}{d\theta} = (I + \frac{1}{\lambda} \nabla_{\phi}^2 \mathcal{L}(\phi, S))^{-1} \quad (۶د)$$

(ب) مشتق مرتبه دوم به دست آمده در الگوریتم MAML به صورت  $\nabla_{\theta}^2 \mathcal{L}(\theta, S)$  می‌باشد این در حالیست که در روش معرفی شده با محاسبه  $\nabla_{\phi}^2 \mathcal{L}(\phi, S)$  مواجه می‌شویم. این بدان معناست که مشتق دوم مورد نظر، کافیت در نقطه بهینه حاصل از حل دقیق بهینه‌سازی ۵ نوشته شود و به مسیری که در گام‌های محاسبه ۵ طی شده است، هیچ وابستگی وجود ندارد. توجه کنید که  $\theta$  همان نقطه شروع بهینه‌سازی است در حالی که  $\phi$  نقطه نهایی بهینه‌سازی محسوب می‌شود. لذا در حل ۵ کافیت تا گرادیان  $\theta$  و  $\phi$  را قطع کنیم، بهینه‌سازی ۵ را حل کرده و گرادیان مورد نیاز برای آپدیت  $\theta$  را با کمک روابط ۳ و ۶ به دست آوریم. این درحالیست که در رویکرد MAML (مخصوصاً زمانی که بیش از یک گام بهینه‌سازی در حلقه داخلی بر می‌داریم)، لازم است تا کل گراف محاسباتی و کل وزن‌های محاسبه شده در میان مسیر را نگه داری کنیم تا بتوانیم مشتق‌های مرتبه دو را محاسبه کنیم.

(ج) • یک وظیفه (Task) به صورت  $\mathcal{T} \in D_{\text{meta-train}}$  نمونه برداری می‌کنیم که شامل داده‌های پشتیبان  $S$  (Support) و پرسمان  $Q$  (Query) می‌باشد.

• با در دست داشتن داده‌های  $S$ ، بهینه‌سازی ۵ را با روش‌های مرسوم بهینه‌سازی انجام می‌دهیم و حاصل را  $\phi$  می‌نامیم.

• طبق رابطه ۶د مشتق  $\phi$  نسبت به  $\theta$  را محاسبه می‌کنیم.

• با کمک داده‌های  $Q$ ، عبارت  $\nabla_{\phi} \mathcal{L}(\phi, Q)|_{\phi=\text{Alg}(\theta, S)}$  را محاسبه می‌کنیم.

• از ضرب دو رابطه محاسبه شده در دو مورد اخیر،  $\nabla_{\theta} \mathcal{L}(\text{Alg}(\theta, S), Q)$  را محاسبه کرده و با کمک این گرادیان متا-پارامترهای  $\theta$  را آپدیت می‌کنیم.

<sup>12</sup>Iterative

روش‌های یادگیری متریک یکی دیگر از روش‌هایی هستند که در درس به عنوان یکی از اعضای خانواده متا-یادگیری با آن‌ها آشنا شدید. در این روش‌ها هدف آن است که یک شبکه استخراج ویژگی مثل  $f_\theta(x)$  یاد گرفته شود تا داده‌ها با برچسب یکسان را در فضای نمایش مخفی در کنار یکدیگر تصویر کند. پارامترهای  $\theta$  متا-پارامترهای مدل در نظر گرفته شده و در طول متا-یادگیری آموزش داده می‌شوند. از آنجایی که  $f_\theta(x)$  صرفاً یک فضای نمایشی<sup>۱۳</sup> فراهم می‌کند، برای کامل کردن شبکه نیاز به یک دسته‌بند (یادگیر پایه<sup>۱۴</sup>) پارامتریک یا نان-پارامتریک داریم که روی این فضای نمایشی قرار گرفته، از داده‌های پشتیبان برای آماده‌سازی خود استفاده کرده و با کمک آن‌ها عمل دسته‌بندی نمونه‌های پرسیمان را انجام دهد. پارامترهای دسته‌بند معرفی شده را به عنوان پارامترهای سریع می‌شناسند و آن‌ها را با  $\phi$  نمایش می‌دهند. این پارامترها مختص هر وظیفه به دست آمده و برای وظیفه بعدی تغییر می‌کنند.

در این روش‌ها یکی از تصمیمات مهم در زمینه طراحی الگوریتم انتخاب مناسب همین دسته‌بند می‌باشد. از جمله انتخاب‌های موجود برای این خانواده، انتخاب روش Nearest Neighbour می‌باشد. همچنین روش دیگری که در مقاله ProtoNet معرفی شد استفاده از دسته‌بندهای مبتنی بر پروتوتایپ دسته‌ها (با میانگین گیری از داده‌های موجود در مجموعه پشتیبان از هر کلاس) می‌باشد. مشکلی که متا-یادگیری مبتنی بر پروتوتایپ دارند این است که دسته‌بند ساده‌ای دارند و تعمیم‌پذیری دسته‌بند در فاز متا-ارزیابی<sup>۱۵</sup> کم می‌باشد. به همین دلیل در مقاله Bertinetto 2018 به بررسی دو دسته‌بند رایج در یادگیری ماشین (Logistic Regression, Ridge Regression) و نحوه استفاده موثر آن در مسئله متا-یادگیری پرداخته شده است. با مطالعه مقاله و راهنمایی‌های داده شده در زیر، به سوالات پاسخ دهید:

(آ) توضیح دهید که در نگاه اول، استفاده از این دسته‌بندها چه مشکلی می‌تواند برای فرایند متا-آموزش<sup>۱۶</sup> ایجاد کند؟ چرا استفاده از رویکردهای پروتوتایپی یا KNN این مشکل را ایجاد نمی‌کند؟ (راهنمایی: پاسخ این مورد غیر مرتبط با پرسش قبل نیست) (۱۰ نمره)

(ب) توضیح دهید که در مقاله معرفی شده، چگونه مشکل معرفی شده را حل می‌کند؟ تفاوت رویکردی که برای Ridge Regression و Logistic Regression به کار گرفته می‌شود را توضیح دهید. (۵ نمره)

(ج) دو ماتریس  $X \in \mathbb{R}^{n \times d}$  به عنوان داده‌های پشتیبان تعبیه شده در فضای نمایش و  $Y \in \mathbb{R}^{n \times l}$  به عنوان برچسب‌های One-hot همان داده‌ها 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 و نحوه استفاده آن در دسته‌بند Ridge Regression را توضیح دهید: (۱۰ نمره)

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

(د) توضیح دهید که برای انجام متا-یادگیری استفاده از کدام یک از دو رابطه بالا بهتر می‌باشد و چرا؟ (۵ نمره)

(ه) در مقاله معرفی شده، برای محاسبه وزن‌های دسته‌بند Logistic Regression از بهینه‌سازی با روش Newton 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. (۱۰ نمره)

(و) تعاریف زیر را در نظر بگیرید:

$$\begin{aligned} A_t &= \text{diag}(q_t) \\ q_t^{(i)} &= \sigma(\omega_t^T x^{(i)}) (1 - \sigma(\omega_t^T x^{(i)})) \\ B_t^{(i)} &= \sigma(\omega_t^T x^{(i)}) - y^{(i)} \end{aligned}$$

که  $A \in \mathbb{R}^{n \times n}$  یک ماتریس قطری و  $B \in \mathbb{R}^n$  می‌باشد.

با اعمال Newton's Method روی تابع هزینه این دسته‌بند، به رابطه به‌روزرسانی زیر برسید (۱۵ نمره):

$$\omega_{t+1} = (X^T A_t X + \lambda I)^{-1} X^T (A_t X \omega_t - B_t)$$

(رابطه فوق مشابه رابطه ۷ مقاله می‌باشد با این تفاوت که به نظر می‌رسد روابط موجود در مقاله در برخی جزئیات از نظر Notation ایراد دارند.)

پاسخ:

(آ) دو دسته بند ذکر شده، دسته‌بندهایی پارامتریک هستند و برای آن که پارامترهای آن‌ها محاسبه شود، نیاز است تا یک دستگاه بهینه‌سازی حل شود. یکی از روش‌های حل دستگاه، استفاده از رویکردهای iterative و استفاده مستقیم از رویکرد Gradient Decent است. بدی استفاده از این نوع Solver آن است که مراحل انجام شده در رویکرد iterative قابلیت عبور گرادیان را ندارند و این باعث به مشکل خوردن متا-یادگیری می‌شود. به عبارت دقیق‌تر، فرض کنید که نمونه‌های پشتیبان را از شبکه استخراج ویژگی عبور داده‌اید و آن‌ها را در فضای

<sup>13</sup>Representation Space

<sup>14</sup>Base-Learner

<sup>15</sup>Meta-Test

<sup>16</sup>Meta-Training

نمایشی تعبیه کرده‌اید. سپس با یک رویکرد iterative بهینه‌سازی مربوط به Ridge Regression را انجام داده و وزن‌های دسته‌بند را به دست می‌آورید. مشکلی که وجود دارد آن است که نمی‌توان گرادیان Meta-Loss را از وزن‌های به دست آمده عبور داد و از طریق آن شبکه استخراج ویژگی و پارامترهای  $\theta$  را آپدیت کرد. این در حالیست که با استفاده از KNN به راحتی می‌توان وزن‌های  $\theta$  را آپدیت کرد چرا که دسته‌بند مستقیماً با خود نمونه‌های هر کلاس کار می‌کند. با استفاده از این رویکرد می‌توان گرادیان را از نمونه‌های موجود هر کلاس عبور داد و از این طریق چینش نمونه‌های یک کلاس را طوری تنظیم کرد که نمونه‌های یک کلاس در فضای نمایش نزدیک یک دیگر بیفتند. رویکرد مشابهی در مورد دسته‌بند پروتوتایپی وجود دارد چرا که عملیات میانگین‌گیری روی نمونه‌های یک کلاس مشکلی برای عبور مسیر گرادیان ایجاد نمی‌کند و با آپدیت کردن  $\theta$  می‌توان نمونه‌ها را طوری در فضا تعبیه کرد که نمونه‌های مربوط به یک کلاس خاص در نزدیکی پروتوتایپ کلاس خودشان قرار بگیرند.

(ب) این مقاله برای حل مشکل Ridge Regression، به جای حل iterative آن، یک جواب Closed-form به عنوان وزن‌های بهینه معرفی می‌کند. این جواب بسته به راحتی قابلیت عبور گرادیان به لایه‌های عقبی را فراهم می‌کند.

از سویی برای دسته‌بند Logistic Regression، از رویکرد iterative خاصی استفاده می‌کند و همزمان نشان می‌دهد که این رویکرد قابلیت عبور گرادیان از تک تک مراحل iterative را فراهم می‌کند. رویکرد مورد استفاده Newton's Method نام دارد.

(ج) 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 هر دو ماتریس‌های مثبت نیمه‌معین می‌باشند و بنابراین حتماً معکوس‌پذیر می‌باشند.

رابطه معرفی شده، همان جواب Closed-Form برای دسته‌بند Ridge Regression است و همانطور که مشاهده می‌شود، در محاسبه این وزن‌ها تنها از ضرب ماتریسی و معکوس‌گیری استفاده شده است که این اپراتورها قابلیت عبور گرادیان از خود را فراهم می‌کنند.

(د) برای انجام متا-یادگیری 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 متا-یادگیری می‌باشد چرا که که تعداد نمونه‌های پشتیبان در مسائل few-shot خیلی کم می‌باشد.

(ه) به دلیل محدودیت تعداد به‌روزرسانی پارامترها در حلقه درونی الگوریتم، از Newton's Method برای همگرایی سریع‌تر به نسبت گرادیان گیری ساده استفاده شده است. در Newton's Method علاوه بر گرادیان، اطلاعات مشتق دوم نیز استفاده می‌شود. توجه کنید که در تابع هزینه‌های محدب (از جمله همین Logistic Regression) استفاده از این رویکرد می‌تواند خیلی سریع ما را به نقطه بهینه همگرا کند. در این روش با نوشتن بسط Taylor برای تابع هزینه  $f$  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 به روزرسانی و حرکت می باشد مشتق میگیریم که جهت بهینه به روزرسانی را پیدا کنیم:

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

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

بنابراین رابطه بهینه سازی را می توان به صورت زیر نوشت:

$$x_{new} = x_{old} - \frac{f'(x_{old})}{f''(x_{old})} = x_{old} - (\nabla_x^2 f)^{-1} \nabla_x f$$

(و) هدف سوال کمینه کردن منفی لگاریتم بیشینه درست نمایی به همراه یک جمله منظم سازی می باشد.

$$P(Y|\omega, X) = \prod_{i=1}^N \left( \sigma(\omega^T x^{(i)}) \right)^{y^{(i)}} \left( 1 - \sigma(\omega^T x^{(i)}) \right)^{1-y^{(i)}}$$

$$L = -\log P(Y|\omega, X) + \lambda \|\omega\|^2 = -\sum_{i=1}^N \left[ y^{(i)} \log \left( \sigma(\omega^T x^{(i)}) \right) + (1 - y^{(i)}) \log \left( 1 - \sigma(\omega^T x^{(i)}) \right) \right] + \lambda \|\omega\|^2$$

حال با گرادیان گیری مقادیر گرادیان اول و دوم را حساب می کنیم:

$$\nabla_{\omega} L = \sum_{i=1}^N \left( \sigma(\omega^T x^{(i)}) - y^{(i)} \right) x^{(i)T}$$

$$H = \nabla_{\omega}^2 L = \sum_{i=1}^N \sigma(\omega^T x^{(i)}) \left( 1 - \sigma(\omega^T x^{(i)}) \right) x^{(i)T} x^{(i)}$$

حال اگر سیگماهای بالا را به ضرب ماتریسی تبدیل کنیم، طبق رابطه به روزرسانی داریم:

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

### سوال ۳: (نظری) استفاده از دسته بند SVM در رویکرد یادگیری متریک (۳۰ نمره)

در این سوال قصد داریم مقاله **MetaOptNet** را مورد بررسی قرار دهیم که برای مدت قابل توجهی SOTA<sup>۱۷</sup> در زمینه یادگیری با تعداد نمونه کم<sup>۱۸</sup> به شمار می رفت. این مقاله شباهت بسیار زیادی به مقاله معرفی شده در پرسش قبل دارد با این تفاوت که قصد دارد به جای دسته بندهای معرفی شده، دسته بند SVM را به عنوان لایه آخر روی شبکه استخراج ویژگی سوار کند. رویکردهایی که در پرسش قبل برای Ridge Regression و Logistic Regression معرفی شدند را نمی توان برای SVM به کار برد. لذا این مقاله به دنبال ارائه روشی است تا این مشکل را برطرف نماید. برای این منظور، رویکردی شبیه به پرسش اول را دنبال می کند. با این تفاوت که در پرسش اول کل وزن های شبکه (اعم از دسته بند و Backbone) در دستگاه بهینه سازی قرار می گرفتند اما در این سوال دستگاه بهینه سازی فقط روی وزن های دسته بند SVM نوشته می شود. برای درک بهتر روش، مقاله مورد نظر را بررسی نموده و به پرسش های زیر پاسخ دهید:

(آ) توضیح دهید که پارامترهای سریع چگونه به کمک داده های پشتیبان و پارامترهای آهسته ساخته می شوند و دستگاه بهینه سازی معرفی شده در این مقاله که از حل آن پارامترهای سریع ساخته می شوند را به همراه دوگان آن به صورت دقیق و با ذکر جزئیات نمادگذاری معرفی کنید. (۵ نمره)

<sup>17</sup>State-of-the-Art

<sup>18</sup>Few-Shot Learning

(ب) با مطالعه صفحه ۴ این مقاله، توضیح دهید که استفاده از قضایای KKT و Implicit Function Theorem چه کمکی در راستای محاسبه گرادیان می‌کند و آپدیت شبکه Backbone با استفاده از چه گرادیانی انجام می‌شود؟ (برای این قسمت اثبات دقیق ریاضی مد نظر نیست و به شرطی که به صورت شفاف کاربرد این دو قضیه و نحوه استفاده آن‌ها را بیان کنید نمره کامل را دریافت می‌کنید) (۱۰ نمره)

(ج) با مطالعه روابط این مقاله ضمن نوشتن دستگاه بهینه‌سازی مربوطه، توضیح دهید که چگونه می‌توان دسته‌بند Ridge Regression مطرح شده در سوال قبل را ذیل همین رویکرد جای داد. (۵ نمره)

(د) بعد از به دست آوردن وزن‌های بهینه  $w$  برای دسته‌بند SVM یا Ridge Regression، تابع متا-هزینه به چه صورتی نوشته می‌شود؟ از این تابع هزینه برای به‌روزرسانی کدام پارامترها استفاده می‌شود؟ (از رابطه ۱۲ مقاله کمک بگیرید، ولی جزئیات به دست آوردن آن را به صورت شفاف بیان کنید) (۱۰ نمره)

توجه: در ادبیات این مقاله، نمادگذاری رایج متا-یادگیری رعایت نشده است و جای نماد  $\theta$  و  $\phi$  با هم عوض شده است. برای یکسان شدن جواب‌ها، شما از نمادگذاری معرفی شده در پرسش اول استفاده کرده و  $\theta$  و  $\phi$  را به ترتیب برای متا-پارامترها و پارامترهای مختص وظیفه به کار ببرید.

پاسخ

(آ) در این مقاله، برای به دست آوردن پارامترهای سریع، پیشنهاد شده است تا یک دستگاه بهینه‌سازی مطابق با تابع هزینه SVM حل شود. در مقاله مورد نظر، رابطه (۴) این دستگاه را نشان می‌دهد که در آن  $w$  ها پارامترهای سریع هستند (لازم است این دستگاه نوشته شود و متغیرهای موجود در آن معرفی شوند). در این رویکرد، پارامترهای آهسته همان وزن‌های شبکه استخراج ویژگی هستند که در این مقاله، بر خلاف نمادگذاری رایج، با  $\phi$  نشان داده شده است. به عبارت دیگر، شبکه استخراج ویژگی یک فضای نمایشی ایجاد می‌کند که دسته‌بند SVM بتواند روی آن قرار گرفته و کلاس‌ها را دسته‌بندی نماید.

(ب) روابط (۴) و (۵) مقاله، دستگاه‌های بهینه‌سازی را نشان می‌دهد که مربوط به بهینه‌سازی Inner Loop می‌باشد. حل این دستگاه‌ها با رویکردهای Iterative شدنی است اما مشکلی که به وجود می‌آید آن است که نمی‌توان به سادگی گرادیان را از مراحل Iterative این دستگاه‌ها به عقب برگرداند تا متا-پارامترها آپدیت شوند.

برای این منظور، به جای آن که گرادیان را از مراحل میانی محاسباتی به صورت عقب گرد عبور دهند، سعی می‌شود تا گرادیان را به صورت مستقیم محاسبه کنند. (پاسخ این سوال شباهت زیادی به پرسش اول دارد) به عبارت دیگر، فرض کنید دستگاه‌های ذکر شده حل شده‌اند و پارامترهای سریع نهایی به دست آمده‌اند. در این صورت، طبق قضیه KKT می‌توان ادعا کرد که مشتق تابع هزینه دستگاه در نقطه بهینه صفر است. از سوی دیگر، اگر از قضیه مشتق ضمنی استفاده شود، می‌توان گرادیان پارامترهای سریع نسبت به متغیرهای نهان فضای نمایش و سپس نسبت به متا-پارامترهای مدل را محاسبه کرد. (رابطه (۸) مقاله)

(ج) اگر Ridge Regression را به صورت یک دستگاه بهینه‌سازی بنویسیم، می‌توان حتی علی‌رغم حل Iterative آن، از رویکرد معرفی شده در این سوال استفاده کرد و محاسبه گرادیان‌ها را ممکن ساخت. (رجوع شود به رابطه ۱۱ مقاله).

(د) با فرض این که پارامترهای سریع  $w$  در مراحل قبلی محاسبه شده‌اند، از تابع هزینه log-likelihood در کنار پارامتر قابل یادگیری  $\gamma$  به عنوان Temperature استفاده می‌شود:

$$\begin{aligned}\mathcal{L}(D^{test}, \theta, \phi, \lambda) &= - \sum_{(x,y) \in D^{test}} \log p(Y=y|x) = - \sum_{(x,y) \in D^{test}} \text{softmax}(\gamma w f_{\theta}(x)) \\ &= - \sum_{(x,y) \in D^{test}} \log \frac{\exp(\gamma w_y f_{\theta}(x))}{\sum_k \exp(\gamma w_k f_{\theta}(x))} = \sum_{(x,y) \in D^{test}} (-\gamma w_y f_{\theta}(x) + \log \sum_k \gamma \exp(w_k f_{\theta}(x)))\end{aligned}$$

از رابطه ذکر شده برای آپدیت متا-پارامترها ( $\theta$ ) استفاده می‌شود. اما باید توجه داشت که روابط فوق علاوه بر وابستگی مستقیم، از طریق  $w$  نیز به  $\theta$  وابستگی دارند. لذا در محاسبه مشتق‌های جزئی، رابطه محاسبه شده در رابطه (۸) مقاله مورد استفاده واقع می‌شود.

#### سوال ۴: (نظری) تنظیم توزیع برای یادگیری چندنمونه‌ای (۱۵ نمره)

یکی از ریسک‌های احتمالی در یادگیری چندنمونه‌ای احتمال بیش برآزش<sup>۱۹</sup> بر روی دادگان کم‌تعداد آموزشی است. در این مقاله روشی پیشنهاد شده است تا به کمک استخراج مشخصات آماری کلاس‌های حاضر در متا-آموزش بتوان توزیع دادگان کلاس‌های حاضر در متا-ارزیابی 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 متا-آموزش 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 متا-آموزش اتخاذ کرده است و چگونه این روش موجب حل مشکل کشیدگی توزیع می‌شود؟ (۵ نمره)

<sup>19</sup>Overfitting



4/1/ در این مقام از **Tukey's ladder of powers transformation** برای این توزیع‌ها بیشتر شبیه به توزیع گوسی شده است. **ساده است**. این تبدیل به صورت زیر تعریف می‌شود.

$$x_\lambda = \begin{cases} x^\lambda & \lambda \neq 0 \\ \log(x) & \lambda = 0 \end{cases}$$

که در آن  $\lambda$  پارت‌هاست. اگر  $\lambda = 1$  به یک تقصیری صورت نمی‌گیرد. کوچک شدن  $\lambda$  باعث می‌شود توزیع کمتر کشیده مثبت (positive skewness) داشته باشد.

(ب) پس از استخراج میانگین و کواریانس کلاس‌های حاضر در متا-آموزش، مدل ارائه شده اقدام به تنظیم توزیع دادگان حاضر در متا-ارزیابی 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 متا-آموزش به کلاس منظور در متا-ارزیابی چه کمکی به تنظیم توزیع می‌کند؟ (۵ نمره)

4/2/ فرض کنید در وضعیت یادگیری برای کلاس‌های مختلف  $\mu_i$  به دست آمده باشد (پس از کس کردن نویسی تبدیل **Tukey ladder of powers**) حال برای داده‌ی  $x$  در درجه‌های مختلف است  $\lambda$  کلاس‌های را که در میان کلاس‌ها (base classes) بیشتر به آن شبیه‌تر باشد پیدا کنیم.

$$S_b = \left\{ -\|\mu_i - \tilde{x}\|_2^2 \mid i \in C_b \right\}$$

$$S_N = \left\{ i \mid -\|\mu_i - \tilde{x}\|_2^2 \in \text{top}_k(S_b) \right\}$$

حال برای این داده‌ی  $x$  به توزیع  $\mu$  کالبرته **calibrated** با این اگورها دست می‌آوریم:

$$\mu' = \frac{\sum_{i \in S_N} \mu_i + \tilde{x}}{k+1}, \quad \sum' = \frac{\sum_{i \in S_N} z_i}{k} + \alpha$$

این تبدیل به یک توزیع  $\mu'$  که در هر چند از این توزیع **calibration** داده‌ها به آن نزدیک‌تر می‌شود.

حال از آنجا که در بازه‌های تعداد داده‌ها است می‌تواند در **inner loop** به راحتی به آن‌ها شبیه‌تر باشد.

برای همین از این توزیع به دست آمده استفاده می‌شود و تعدادی از این توزیع‌ها برای می‌شود و در آن‌ها استفاده می‌شود.

همچنین برای توضیح بیشتر فرض کنید در هنگام متاست تکلیف ما جدا کردن کلاس روباه از کلاغ در تنظیمات نمونه کم باشد. در صورتی که ما در متاآموزش کلاس‌های گرگ و کبوتر را داشته باشیم نمونه‌های کلاس‌های گرگ و کبوتر به ترتیب می‌توانند به کلاس‌های روباه و کلاغ به علت شباهت کلاس‌ها کمک کنند و مانع از آورفیت شدن شبکه بر روی تعداد کم داده‌ها شوند.

(ج) در تنظیماتی که در هنگام متا-ارزیابی 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 تنظیم توزیع جداگانه چه مزیتی نسبت به میانگین گیری نمونه‌ها و سپس یک تنظیم توزیع دارد؟ (۵ نمره)

پاسخ

ج۴) این کار، در مزیت اصلی دارد. اولاً یک سمپل می‌تواند کل توزیع را بایاس کند، به عبارت دیگر آن میانگین بگیریم ممکن است میانگین به سبک سمپل غامس بایاس شود. زمانی که با تنظیم جداگانه، این اتفاق رخ نمی‌دهد. ثانیاً به دلیل تنوع در توزیع که هر سمپل از یک کلاس می‌دهد، سمپل می‌کند که در فاز تست می‌گیریم، متنوع تر می‌شود و بنابراین هم بهتر و متنوع تر می‌باشی آن کلاس فراهم داشت.

به علاوه در صورتی که توزیع دادگان کلاس حاضر در متاتست یک توزیع چند قله ای باشد، استفاده از یک نمونه نمی‌تواند مدل خوبی از آن توزیع چند قله ای ارائه دهد. در حالی که استفاده از چند نمونه این مشکل را حل می‌کند.

#### سوال ۵: (عملی) یادگیری چند نمونه‌ای از طریق یادگیری متریک (۲۵ 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 نماییم. سپس این دادگان را به سه زیر مجموعه متا-آموزش، متا-اعتبارسنجی<sup>۲۰</sup> و متا-ارزیابی 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 Way و Shot بتواند دادگان اتکا و پرسمان برای یک وظیفه را تولید کنند (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 Backbone استخراج ویژگی استفاده کنید و در جلوی آن دو لایه تمام متصل<sup>۲۱</sup> با اندازه دلخواه قرار دهید.

(آ) دسته‌بند را با تنظیمات 8-shot, 10-way آموزش دهید و سپس دقت مدل را بر روی دادگان متا-ارزیابی گزارش دهید. انتظار می‌رود دقت در این بخش بیشتر از ۵۱ درصد باشد.

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

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

(د) حال در هنگام متا-آموزش با تنظیمات 10-shot, 10-way دسته‌بند را آموزش دهید. در هنگام متا-ارزیابی اما متا-ارزیابی 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 ۶: (عملی) متا-یادگیری براساس بهینه‌سازی (۵۰ نمره)

در این سوال قصد داریم تا مدل معروف دسته متا-یادگیری 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 متا-یادگیری 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 متا-یادگیری بخش اول + ساختار مدل pretrained داده شده) را به صورت متاپارامتر در نظر بگیرید و متا-یادگیری را روی آن انجام دهید. در نهایت نتایج بدست آمده از هر دو بخش را تحلیل و گزارش نمایید.

<sup>20</sup>Meta-Validation

<sup>21</sup>Fully Connected