

بسمه تعالی

پاسخ سری ششم تمرینات درس یادگیری ماشین

امیرحسین رمضانی بناب (۹۹۲۱۰۲۹۴)

۱ سوال ۱

۱.۱ آ

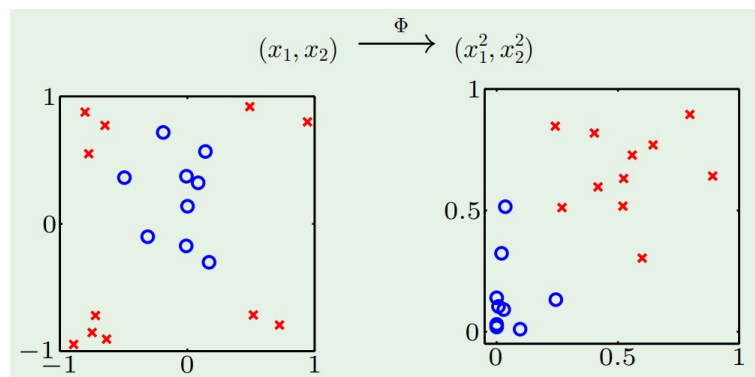
۱. اصل *Occam's Razor*: این اصل بیان می‌کند که ساده‌ترین مدل برای مدل کردن یک مجموعه‌ی داده، محتمل‌ترین مدل است. یعنی مدل‌هایی که H ساده‌تری دارند، دقت‌های *Out Of Sample* مناسبی خواهند داشت. شهود این اصل هم به این صورت است که مدل‌های ساده‌تر به دلیل اینکه تعداد کمتری نسبت به مدل‌های پیچیده‌تر دارند، احتمال کمتری دارد که به یک داده فیت شوند. چون فیت شدن به داده سخت‌تر است، پس موقعی که یک مدل ساده به یک مجموعه‌ی داده فیت می‌شود ارزش بالاتری دارد و نسبت به مدل پیچیده ارجح است.

۲. اصل *Sampling Bias*: طبق این اصل، داده‌ها باید به شکل عادلانه‌ای جمع‌آوری شده باشند و نمونه‌برداری آن‌ها به درستی انجام شده باشد. به عنوان مثال اگر جمع‌آوری داده‌ها به شکل *biased* انجام شده باشد، موقع عملکرد مدل در دنیای واقعی، داده‌ها از توزیع دیگری خواهند بود که با توزیع داده‌های آموزشی تفاوت دارد. این تفاوت منجر به شکست مدل در دنیای واقعی مسئله می‌شود.

۳. اصل *Data Snooping*: طبق این اصل، اگر دیتاست به هر نحوی روی هر مرحله‌ای از فرایند یادگیری تاثیر گذاشته باشد، ارزیابی خروجی الگوریتم روی این دیتاست باعث تولید دقت‌های ناعادلانه می‌گردد. زیرا تاثیر پذیرفتن بدون حسابرسی بخش‌هایی از فرایند یادگیری از روی داده به این معناست که بخشی از فضای فرضیه‌ی مسئله با دیدن این داده‌ها کنار گذاشته می‌شود و در نهایت مدل دقیقی به دست می‌آید که فقط روی داده‌های آموزشی خوب عمل می‌کند. در حالی که یک مدل مناسب باید بتواند قابلیت تعمیم‌پذیری داشته باشد.

۲.۱ ب

خیر. فرض کنید تبدیلی انجام داده‌ایم که داده‌ها را از فضای \mathcal{X} به فضای \mathcal{Z} برده است. ممکن است داده‌ها در فضای \mathcal{X} دارای مزر تصمیم‌گیری بسیار پیچیده‌ای با یک الگوریتم خطی باشند ولی در فضای \mathcal{Z} از هم خطی جداپذیر باشند. به عنوان مثال شکل زیر که از اسلاید درس وام گرفته شده است را آورده‌ام.



همانطور که مشخص است مرز تصمیم‌گیری تصویر سمت چپ یک مرز پیچیده خواهد بود. در حالی که با تصویر کردن داده‌ها در فضای جدید به یک مرز ساده‌ی خطی خواهیم رسید. پس اینکه یک مرز پیچیده داشته باشیم لزوماً به این معنا نیست که اصل *Occam's Razer* نقض می‌شود.

۳.۱ ب

(ایده گرفته شده از کتاب آقای ابومصطفی)

از آنجایی که *Sampling Bias* مربوط به جمع‌آوری داده و *Data Snooping* مربوط به انتخاب مدل یادگیری است، این دو مفهوم با هم تفاوت دارند و به فازهای متفاوتی از آموزش ارتباط دارند. ولی مواقعی پیش می‌آید که بین این دو مفهوم، ارتباط برقرار می‌شود. مثلاً فرض کنید یک بانک فقط با یک مجموعه از شرکت‌ها سر و کار دارد که هر سال به هر یک از این شرکت‌ها میزان مشخصی وام تخصیص می‌دهد. بانک می‌خواهد برای برنامه‌ریزی میزان وام‌های سال جاری، از یادگیری ماشین کمک بگیرد. بدین منظور داده‌های ۲۰ سال اخیر شرکت‌هایی که در حال حاضر با بانک کار می‌کنند را در نظر می‌گیرد و میزان وامی که به آن شرکت‌ها تخصیص خواهد داد را پیش‌بینی می‌کند. اما این الگوریتم به احتمال قوی دقت مناسبی نخواهد داشت. زیرا شرکت‌هایی که بانک در ۲۰ سال اخیر با آن‌ها قطع رابطه کرده‌است در نظر گرفته نشده‌اند و تنها شرکت‌هایی لحاظ شده‌اند که در ۲۰ سال اخیر با بانک در ارتباط بوده‌اند. این مشکل همان *Sample Bias* است. اگر به ریشه‌ی آن توجه کنیم، متوجه می‌شویم که نباید داده‌های بقیه‌ی شرکت‌ها را حذف می‌کردیم و با گرفتن زیرمجموعه‌ای از داده شرکت‌ها از روی تاریخچه‌ی آن‌ها به نوعی دچار مشکل *Data Snooping* شده‌ایم. پس در مواقعی *Data Snooping* می‌تواند منجر به *Sample Bias* شود.

۲ سوال ۲

ابتدا مسئله‌ی *Locally Weighted Linear Regression* را بیان می‌کنیم. تنها تفاوت این مسئله با مسئله‌ی رگرسیون خطی ساده این است که معیار خطای استفاده شده در آن به شکل زیر است

$$J(\theta) = \frac{1}{2} \sum_{n=1}^N w_n (\theta^T \mathbf{x}_n - y_n)^2 = \frac{1}{2} (\mathbf{X}\theta - Y)^T \mathbf{W} (\mathbf{X}\theta - Y)$$

که w_n یک عدد نامنفی است که میزان تاثیری که برچسب داده‌ایی که می‌خواهیم پیش‌بینی کنیم از برچسب داده‌ی آموزشی n ام می‌پذیرد را نشان می‌دهد. \mathbf{W} ماتریس قطری $n \times n$ است که

$$\mathbf{W}[i, i] = w_i$$

همچنین چون ماتریس قطری است پس

$$\mathbf{W}^T = \mathbf{W}$$

۱.۲ آ

برای به دست آوردن θ بهینه، از J نسبت به θ مشتق می‌گیریم و برابر صفر قرار می‌دهیم.

$$\begin{aligned} \frac{\partial J}{\partial \theta} &= \frac{1}{2} \frac{\partial}{\partial \theta} (\mathbf{X}\theta - Y)^T \mathbf{W} (\mathbf{X}\theta - Y) \\ &= \frac{1}{2} \frac{\partial}{\partial \theta} [(\theta^T \mathbf{X}^T - Y^T) \mathbf{W} (\mathbf{X}\theta - Y)] \\ &= \frac{1}{2} \frac{\partial}{\partial \theta} [\theta^T \mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \theta^T \mathbf{X}^T \mathbf{W} Y - Y^T \mathbf{W} \mathbf{X} \theta + Y^T \mathbf{W} Y] \\ &= \frac{1}{2} \frac{\partial}{\partial \theta} [\theta^T \mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \theta^T \mathbf{X}^T \mathbf{W} Y - Y^T \mathbf{W} \mathbf{X} \theta] \end{aligned}$$

حال اگر عبارت بالا را برابر صفر قرار دهیم داریم:

$$\begin{aligned} \frac{1}{2} \frac{\partial}{\partial \theta} [\theta^T \mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \theta^T \mathbf{X}^T \mathbf{W} Y - Y^T \mathbf{W} \mathbf{X} \theta] &= 0 \\ \Rightarrow \frac{1}{2} [2\mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \mathbf{X}^T \mathbf{W} Y - \mathbf{X}^T \mathbf{W}^T Y] &= 0 \\ \Rightarrow \frac{1}{2} [2\mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \mathbf{X}^T \mathbf{W} Y - \mathbf{X}^T \mathbf{W} Y] &= 0 \\ \Rightarrow \frac{1}{2} [2\mathbf{X}^T \mathbf{W} \mathbf{X} \theta - 2\mathbf{X}^T \mathbf{W} Y] &= 0 \\ \Rightarrow \mathbf{X}^T \mathbf{W} \mathbf{X} \theta - \mathbf{X}^T \mathbf{W} Y &= 0 \\ \Rightarrow \mathbf{X}^T \mathbf{W} \mathbf{X} \theta &= \mathbf{X}^T \mathbf{W} Y \end{aligned}$$

حال برای به دست آوردن مقدار θ دو طرف معادله‌ی بالا را [با فرض معکوس پذیر بودن] از سمت چپ در $(\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$ ضرب می‌کنیم. داریم

$$\theta = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} Y$$

۱. وزن‌های W باعث می‌شود که برای پیش‌بینی برچسب یک داده، از همه‌ی داده‌های آموزشی به شکل یکسان بهره نبریم. به این شکل که هر داده‌ی آموزشی، براساس شباهتی که با داده‌ی آزمون دارد، در انتخاب برچسب برای آن دخالت می‌کند. برای محاسبه‌ی شباهت نیز می‌توان از معیارهای متفاوتی استفاده کرد و طبق معیار انتخاب شده، شباهت داده‌ی آزمون با داده‌های آموزش در W قرار داد.
۲. به واریانس W گوسی، پهنای باند می‌گویند زیرا هر چقدر مقدار واریانس بیشتر شود، داده‌های دورتر با وزن بیشتری در انتخاب برچسب داده‌ی جدید دخیل هستند و واریانس کوچک تر باعث می‌شود که داده‌های نزدیک‌تر اهمیت بیشتری پیدا کنند و داده‌های دورتر اهمیت کمتر. پس واریانس، سرعت افت مقادیر W را بر حسب فاصله از داده‌ی آزمون نمایش می‌دهد. داده‌هایی که وزن بیشتری بگیرند در تعیین fit تاثیر بیشتری خواهند گذاشت.
۳. واریانس این توزیع می‌تواند مقادیر متمایزی برای هر داده‌ی آزمون داشته باشد. این متغیر کردن توزیع می‌تواند به طرق مختلف انجام شود. به عنوان مثال می‌توان ابتدا یک الگوریتم خوشه بندی روی داده‌های آموزش اجرا کرد و سپس واریانس همه‌ی داده‌های آزمونی که به یک خوشه‌ی داده‌های آموزشی تعلق می‌گیرند را یکسان در نظر گرفت و آن مقدار را به نحوی تنظیم کرد که داده‌های آزمون داخل هر خوشه، بیشترین اثر را از داده‌های آموزش داخل آن خوشه بگیرند. به این ترتیب احتمالاً عملکرد الگوریتم بهبود می‌یابد.

برای این بخش ابتدا الگوریتم کلی حالت *Regression* بیان می شود و سپس با اعمال تغییراتی برای *Classification* نیز استفاده می گردد.

۱. حالت *Regression*:

- فرض می شود ورودی مسئله دادگان آموزشی $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ باشند. همچنین یک تابع *loss* مشتق پذیر $L(\mathbf{y}, f(\mathbf{x}))$ داریم.
- ابتدا یک مقدار اولیه ثابت برای پیش بینی انتخاب می کنیم. این مقدار به شکلی انتخاب می شود که معیار *loss* را کمینه کند.

$$F_0(\mathbf{x}) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(\mathbf{y}_i, \gamma)$$

با این کار یک نقطه ی شروع خوب برای الگوریتم داریم.

- حال به سراغ ساخت یادگیرنده های جدید می رویم. اگر در مرحله ی ساخت یادگیرنده ی m ام باشیم این مراحل را پیش می بریم
- ابتدا *residual* ها را محاسبه می کنیم. به این معنا که خطای الگوریتم روی داده های مختلف را اندازه می گیریم.

$$r_{im} = - \left[\frac{\partial L(\mathbf{y}_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}$$

به این شکل *residual* الگوریتم در مرحله ی m روی داده ی i محاسبه می شود.

- حال سعی می کنیم توسط یک یادگیرنده ی پایه، این مقادیر *residual* را یاد بگیریم و با افزودن این یادگیرنده، بتوانیم در راستای کاهش خطا حرکت کنیم. نام تابع به دست آمده از این یادگیرنده را $h_m(\mathbf{x})$ می گذاریم.
- در قدم بعدی سعی می کنیم وزنی که یادگیرنده ی جدید در کل مدل خواهد داشت را انتخاب کنیم. این وزن طوری انتخاب می شود که در راستای کاهش خطای مدل حرکت کند.

$$\gamma_m = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i) + \gamma h_m(\mathbf{x}_i))$$

- در نهایت یادگیرنده ی ساخته شده را با وزن به دست آمده به کل مدل اضافه می کنیم

$$F_{m-1}(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \gamma h_m(\mathbf{x})$$

۲. حالت *Classification*: در این حالت فرض کنید که M کلاس داریم. در آن صورت به جای یک تابع F ، به تعداد M تابع خواهیم داشت که هر یک از این توابع امتیاز حضور یک داده در هر کلاس را مشخص می کنند. یعنی توابع F_1, F_2, \dots, F_M را داریم که $F_i(\mathbf{x}_j)$ نشان دهنده ی امتیاز حضور داده ی j در کلاس i است. سپس از روی این مقادیر امتیاز، می توانیم توسط تابع *Softmax* احتمال حضور داده در هر کلاس را مشخص کرده و داده را به کلاسی تخصیص دهیم که دارای بیشترین احتمال باشد. به این شکل که برای داده ی \mathbf{x} ، احتمال حضور آن در هر کلاس را به شکل زیر محاسبه می کنیم.

$$P_m(\mathbf{x}) = \frac{e^{F_m(\mathbf{x})}}{\sum_{c=1}^M e^{F_c(\mathbf{x})}}, \quad m = 1, \dots, M$$

و برچسب داده برابر $\underset{m \in \{1, \dots, M\}}{\operatorname{argmax}} P_m(\mathbf{x})$ خواهد بود.

مراحل ابتدایی الگوریتم به شکل زیر است.

- ابتدا هر برچسب را به یک توزیع احتمال تبدیل می کنیم. یعنی برای داده ی i ام

$$Y_1(\mathbf{x}_i) = 0, Y_2(\mathbf{x}_i) = 0 \dots Y_{y_i}(\mathbf{x}_i) = 1 \dots Y_M(\mathbf{x}_i) = 0$$

- طبق مقادیر بالا، احتمال عضویت داده در هر کلاس $(P_1(\mathbf{x}_i), \dots, P_M(\mathbf{x}_i))$ را طبق فرمول ارائه شده محاسبه می‌کنیم.
- اختلاف مقادیر پیش‌بینی شده و مقادیر واقعی را به عنوان *residual* در نظر می‌گیریم.

حال به تفاوت مراحل اصلی الگوریتم می‌پردازیم

- در حالت *Regression*، یک بردار F برای داده‌ها داشتیم که می‌خواستیم مقدار آن را بهینه کنیم. در حالت *Classification*، با یک ماتریس طرف هستیم

$$\begin{bmatrix} F_1(\mathbf{x}_1) & F_2(\mathbf{x}_1) & \dots & F_M(\mathbf{x}_1) \\ F_1(\mathbf{x}_2) & F_2(\mathbf{x}_2) & \dots & F_M(\mathbf{x}_2) \\ \dots & \dots & \dots & \dots \\ F_1(\mathbf{x}_n) & F_2(\mathbf{x}_n) & \dots & F_M(\mathbf{x}_n) \end{bmatrix}$$

- در نتیجه برای fit کردن یک مدل روی *residuals* نیز باید از یک ماتریس برای نشان داده‌گرایان‌ها استفاده کنیم.

$$\begin{bmatrix} \frac{\partial L}{\partial F_1(\mathbf{x}_1)} & \frac{\partial L}{\partial F_2(\mathbf{x}_1)} & \dots & \frac{\partial L}{\partial F_M(\mathbf{x}_1)} \\ \frac{\partial L}{\partial F_1(\mathbf{x}_2)} & \frac{\partial L}{\partial F_2(\mathbf{x}_2)} & \dots & \frac{\partial L}{\partial F_M(\mathbf{x}_2)} \\ \dots & \dots & \dots & \dots \\ \frac{\partial L}{\partial F_1(\mathbf{x}_n)} & \frac{\partial L}{\partial F_2(\mathbf{x}_n)} & \dots & \frac{\partial L}{\partial F_M(\mathbf{x}_n)} \end{bmatrix}$$

- پس برای هر *iteration* به جای مقدار γ ، به تعداد کلاس‌ها پارامتر خواهیم داشت تا همه‌ی توابع را در راستای عکس‌گردان حرکت دهیم.
- در نهایت، در پایان هر *iteration* تمام توابع با وزن‌های مربوط به فرضیه‌های جدید بروزرسانی خواهند شد.

$$F_i = F_i + \gamma_i h_i$$

روش‌های دیگری نیز برای فرموله‌بندی مسئله‌ی دسته‌بندی وجود دارد. روشی که در این پاسخ ارائه شد، برگرفته از این لینک است.

۲.۳ ب

در حالت دو کلاسه، به جای استفاده از *One Hot*، می‌توانیم از تابع هزینه‌ی *Cross Entropy* استفاده کنیم. اما به دلیل اینکه ورودی دوم تابع هزینه از جنس احتمال نمی‌باشد، باید تغییراتی در تابع هزینه اعمال شود. در صورتی که فرض کنیم "امتیاز" محاسبه شده تا این لحظه برای تخصیص داده‌ی i ام به کلاس ۱ برابر s_i باشد. در این صورت تابع *sigmoid* را روی این مقدار اعمال کرده و آن را برابر "احتمال" تخصیص داده‌ی i به کلاس ۱ در نظر می‌گیریم و با p_i نمایش می‌دهیم. پس تابع هزینه به شکل زیر می‌باشد

$$L(\mathbf{y}, s) = - \sum_{i=1}^n \mathbf{y}_i \log(\sigma(s_i)) + (1 - \mathbf{y}_i) \log(1 - \sigma(s_i))$$

به عبارتی برای داده‌ی i ام تابع هزینه برابر است با

$$-[\mathbf{y}_i \log(\sigma(s_i)) + (1 - \mathbf{y}_i) \log(1 - \sigma(s_i))]$$

رابطه‌ی بالا را کمی ساده‌تر می‌کنیم.

$$\begin{aligned} -[\mathbf{y}_i \log(\sigma(s_i)) + (1 - \mathbf{y}_i) \log(1 - \sigma(s_i))] &= -\mathbf{y}_i \log(\sigma(s_i)) - (1 - \mathbf{y}_i) \log(1 - \sigma(s_i)) \\ &= -\mathbf{y}_i \log\left(\frac{e^{s_i}}{1 + e^{s_i}}\right) - (1 - \mathbf{y}_i) \log\left(\frac{1}{1 + e^{s_i}}\right) \\ &= -\mathbf{y}_i [\log(e^{s_i}) - \log(1 + e^{s_i})] - (1 - \mathbf{y}_i) [\log(1) - \log(1 + e^{s_i})] \\ &= -\mathbf{y}_i [s_i - \log(1 + e^{s_i})] + (1 - \mathbf{y}_i) [\log(1 + e^{s_i})] \\ &= -\mathbf{y}_i s_i + \log(1 + e^{s_i}) \end{aligned}$$

اگر مشتق این تابع را نسبت به s_i محاسبه کنیم داریم:

$$\begin{aligned}\frac{d}{ds_i}[-\mathbf{y}_i s_i + \log(1 + e^{s_i})] &= -\mathbf{y}_i + \sigma(s_i) \\ &= -\mathbf{y}_i + p_i\end{aligned}$$

و اگر مشتق دوم این تابع را نسبت به s_i محاسبه کنیم داریم:

$$\begin{aligned}\frac{d^2}{ds_i^2}[-\mathbf{y}_i s_i + \log(1 + e^{s_i})] &= \frac{d}{ds_i}[-\mathbf{y}_i + \sigma(s_i)] \\ &= \frac{e^{s_i}}{(1 + e^{s_i})^2} \\ &= \frac{e^{s_i}}{1 + e^{s_i}} \times \frac{1}{1 + e^{s_i}} \\ &= \sigma(s_i)(1 - \sigma(s_i)) \\ &= p_i(1 - p_i)\end{aligned}$$

حال با استفاده از بسط تیلور عبارت داده شده را ساده می کنیم.

$$\begin{aligned}L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i) + \gamma h_m(\mathbf{x}_i)) &\approx L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) + \frac{d}{dF_{m-1}(\mathbf{x}_i)} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) \gamma h_m(\mathbf{x}_i) \\ &\quad + \frac{1}{2} \frac{d^2}{dF_{m-1}(\mathbf{x}_i)^2} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) \gamma^2 h_m(\mathbf{x}_i)^2\end{aligned}$$

پس تابعی که می خواهیم نسبت به γ کمینه کنیم به شکل زیر است.

$$\sum_{i=1}^n \left[L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) + \frac{d}{dF_{m-1}(\mathbf{x}_i)} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) \gamma h_m(\mathbf{x}_i) + \frac{1}{2} \frac{d^2}{dF_{m-1}(\mathbf{x}_i)^2} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) \gamma^2 h_m(\mathbf{x}_i)^2 \right]$$

که اگر کمی ساده سازی کنیم

$$\begin{aligned}\sum_{i=1}^n [L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i))] &+ \gamma \sum_{i=1}^n \left[\frac{d}{dF_{m-1}(\mathbf{x}_i)} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) h_m(\mathbf{x}_i) \right] \\ &+ \frac{1}{2} \gamma^2 \sum_{i=1}^n \left[\frac{d^2}{dF_{m-1}(\mathbf{x}_i)^2} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) h_m(\mathbf{x}_i)^2 \right]\end{aligned}$$

پس اگر نسبت به γ مشتق بگیریم و برابر صفر قرار دهیم داریم:

$$\sum_{i=1}^n \left[\frac{d}{dF_{m-1}(\mathbf{x}_i)} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) h_m(\mathbf{x}_i) \right] + \gamma \sum_{i=1}^n \left[\frac{d^2}{dF_{m-1}(\mathbf{x}_i)^2} L(\mathbf{y}_i, F_{m-1}(\mathbf{x}_i)) h_m(\mathbf{x}_i)^2 \right] = 0$$

حال اگر مشتقات را جایگذاری کنیم داریم

$$\sum_{i=1}^n [(-\mathbf{y}_i + \sigma(F_{m-1}(\mathbf{x}_i))) h_m(\mathbf{x}_i)] + \gamma \sum_{i=1}^n [\sigma(F_{m-1}(\mathbf{x}_i))(1 - \sigma(F_{m-1}(\mathbf{x}_i))) h_m(\mathbf{x}_i)^2] = 0$$

پس

$$\gamma = \frac{-\sum_{i=1}^n [(-\mathbf{y}_i + \sigma(F_{m-1}(\mathbf{x}_i))) h_m(\mathbf{x}_i)]}{\sum_{i=1}^n [\sigma(F_{m-1}(\mathbf{x}_i))(1 - \sigma(F_{m-1}(\mathbf{x}_i))) h_m(\mathbf{x}_i)^2]}$$

۳.۳ پ

مراحل اجرای الگوریتم خواسته شده در قالب اسلایدی در فایل *Q3p.pdf* به پیوست تقدیم شده است.

۴.۳ ت

مراحل اجرای الگوریتم خواسته شده در قالب اسلایدی در فایل *Q3t.pdf* به پیوست تقدیم شده است.

۴ سوال ۴

نوت‌بوک خواسته شده در قالب فایل *ML2021S_HW6* به پیوست تقدیم شده است.