



یادگیری ماشین

بهار ۱۴۰۴

استاد:

آزمون پایانترم

سوال ۱ حل.

۱.۱) ساختار دو درخت تصمیم یکسان خواهد بود. زیرا که درخت تصمیم بر اساس مقادیر ویژگی‌ها و ترتیب نسبی آن‌ها تقسیم‌بندی انجام می‌دهد، نه مقیاس عددی آن‌ها. به عبارت دیگر، الگوریتم درخت تصمیم به دنبال یافتن آستانه‌هایی است که داده‌ها را به خوبی تفکیک کند، و این آستانه‌ها از روی مقایسه‌ی مقادیر استخراج می‌شوند نه مقدار مطلق ویژگی‌ها. از طرفی استانداردسازی ترتیب داده‌ها را حفظ می‌کند. در فرآیند استانداردسازی، فقط میانگین و انحراف معیار تغییر داده می‌شود، اما ترتیب نسبی مقادیر هر ویژگی ثابت می‌ماند. بنابراین، درخت تصمیم همان تقسیم‌بندی را انجام خواهد داد. از طرفی دیگر معیارهای ارزیابی تقسیم (مانند Gini و آنتروپی) به مقیاس عددی وابسته نیستند. در نتیجه این معیارها براساس توزیع نمونه‌ها در گره‌ها محاسبه می‌شوند و نه مقادیر دقیق ویژگی‌ها و بنابراین تغییر در مقیاس داده‌ها اثری بر انتخاب تقسیم‌ها نخواهد داشت.

۳.۱) پیشنهاد او از لحاظ مفهومی دچار مشکل است. یکی از فرضیاتی که در رابطه با دیتای تست داریم این است که ما این داده‌ها و اطلاعات مربوط آن‌ها را در زمان آموزش در اختیار نداریم و تنها در هنگام تست می‌توانیم این داده‌ها را دسترسی داشته باشیم. بنابراین اگر الگوریتم PCA را استفاده کنیم، از اطلاعات موجود در داده‌ی تست برای نرمال‌سازی و یافتن بردارهای اصلی استفاده کرده‌ایم و نشت اطلاعات داریم.

۴.۱) خیر، تغییری ایجاد نمی‌شود. علت این امر این است که ویژگی‌ای که مقدار ثابتی (مانند یک) داشته باشد، دارای واریانس صفر می‌باشد و در الگوریتم PCA ویژگی‌های با بیشترین میزان واریانس انتخاب می‌شوند و در نتیجه این ویژگی تأثیری نخواهد داشت.

۵.۱) درست است چون PCA در واقع L بزرگترین مقدار ویژه را در نظر گرفته و از بردار ویژه‌های نظیر آن‌ها استفاده می‌کند؛ پس با دوبار انجام دادنش، هم‌چنان به همان بردارها با بزرگ‌ترین مقدار ویژه می‌رسیم.

۶.۱) احتمال این که نمونه‌ی i ام در بین نمونه‌های انتخاب شده نباشد برابر با $(1 - \frac{1}{N})^{Np}$ می‌باشد که برای N های بزرگ این مقدار به e^{-p} میل می‌کند. حال اگر متغیر تصادفی X را برابر با تعداد نمونه‌هایی که در نمونه‌گیری قرار ندارند بنامیم، خواهیم داشت:

$$E[X] = E[X_1 + X_2 + \dots + X_N] = NE[X_1] = N.e^{-p}$$

سوال ۲ حل.

۱۴ شهریور ۱۴۰۴

تعریف ۱ (مسئله اولیه SVM با حاشیه سخت) فرض کنید یک مجموعه داده خطی پذیر $\{(x_i, y_i)\}_{i=1}^N$ داشته باشیم که در آن $x_i \in \mathbb{R}^2$ و $y_i \in \{+1, -1\}$ باشد و حداقل یک نقطه از هر کلاس وجود داشته باشد. پارامترهای ابرصفحه بهینه $w \in \mathbb{R}^2$ و $b \in \mathbb{R}$ با حل مسئله زیر به دست می‌آیند:

$$\min_{w, b} \quad \frac{1}{2} \|w\|^2$$

to subject $y_i(w \cdot x_i + b) \geq 1 \quad \text{برای } i = 1, 2, \dots, N.$

یک بردار پشتیبان نقطه‌ای است که قید مربوط به آن فعال باشد، یعنی $y_i(w^* \cdot x_i + b^*) = 1$.

قضیه ۱ برای $N \geq 3$ نقطه در \mathbb{R}^2 ، پاسخ SVM با حاشیه سخت حداقل ۳ بردار پشتیبان دارد.

اثبات در سه بخش ارائه می‌شود: (۱) اثبات یکتایی بردار وزن بهینه w^* ، (۲) اثبات یکتایی بایاس بهینه b^* ، و (۳) نشان دادن اینکه این یکتایی مستلزم وجود حداقل سه بردار پشتیبان است.

بخش ۱ (یکتایی w^*) تابع هدف $f(w) = \frac{1}{2} \|w\|^2$ به شدت محدب است زیرا هسین آن $\nabla^2 f(w) = I$ (ماتریس همانی) مثبت معین است. قیود نیز همگی خطی هستند و بنابراین یک ناحیه شدنی محدب را تعریف می‌کنند. یک تابع به شدت محدب که روی یک مجموعه محدب کمینه شود، یک کمینه‌سراسری یکتا دارد. در نتیجه، w^* یکتا است.

بخش ۲ (یکتایی b^*) فرض کنید برای تناقض، دو جواب بهینه متفاوت وجود داشته باشند: (w^*, b_1) و (w^*, b_2) به طوری که $b_1 \neq b_2$. هر دو باید تمام قیود را ارضا کنند:

$$y_i(w^* \cdot x_i + b_1) \geq 1 \quad \forall i \in \{1, \dots, N\},$$

$$y_i(w^* \cdot x_i + b_2) \geq 1 \quad \forall i \in \{1, \dots, N\}.$$

با تفاضل قید اول از دوم برای هر نقطه دلخواه i داریم:

$$y_i(w^* \cdot x_i + b_1) - y_i(w^* \cdot x_i + b_2) \geq 0,$$

$$y_i(b_1 - b_2) \geq 0.$$

این نامساوی باید برای تمام نقاط برقرار باشد. با تحلیل بر اساس برچسب کلاس:

- برای هر i با $y_i = +1$: $b_1 \geq b_2$ $\implies 1 \cdot (b_1 - b_2) \geq 0$.
- برای هر i با $y_i = -1$: $b_1 \leq b_2$ $\implies -1 \cdot (b_1 - b_2) \geq 0$.

تنها مقداری که بتواند همزمان $b_1 \leq b_2$ و $b_1 \geq b_2$ را برای تمام نقاط (با فرض وجود حداقل یک نقطه از هر کلاس) برآورده کند، $b_1 = b_2$ است. این با فرض اولیه در تناقض است. بنابراین، b^* نیز باید یکتا باشد.

بخش ۳ (لزوم وجود حداقل سه بردار پشتیبان) شرایط KKT برای جواب بهینه (w^*, b^*) بیان می‌کند که:

(آ) w^* را می‌توان به صورت ترکیب خطی بردارهای پشتیبان نوشت: $w^* = \sum_{i=1}^N \alpha_i y_i x_i$ ، که در آن $\alpha_i > 0$ تنها اگر x_i یک بردار پشتیبان باشد.

(ب) برای هر بردار پشتیبان x_j ، تساوی زیر برقرار است: $y_j(w^* \cdot x_j + b^*) = 1$.

معادله مربوط به یک بردار پشتیبان را می‌توان به صورت زیر بازنویسی کرد:

$$y_j b^* + w_1^*(x_j)_1 + w_2^*(x_j)_2 = 1, \quad (1)$$

این یک معادله خطی با سه مجهول است: w_1^* ، w_2^* و b^* .

جواب یکتا (w^*, b^*) کاملاً توسط مجموعه معادلاتی از نوع (۱) که توسط بردارهای پشتیبان ارائه می‌شوند، تعیین می‌گردد. برای تعیین یکتای سه مجهول، یک دستگاه معادلات نیاز به حداقل سه معادله مستقل خطی دارد. از آنجایی که داده‌ها عمومی هستند و نقاطی از هر دو کلاس وجود دارند، حداقل تعداد بردارهای پشتیبان مورد نیاز برای تشکیل یک دستگاه با رتبه کامل، سه عدد است. داشتن کمتر از سه بردار پشتیبان منجر به یک دستگاه معادلات تحت‌تعیین می‌شود که با یکتایی (w^*, b^*) که در بخش‌های ۱ و ۲ اثبات شد، در تناقض است.

بنابراین، جواب SVM با حاشیه سخت باید حداقل سه بردار پشتیبان داشته باشد.

سوال ۳ حل. رابطه‌ی بین این ۳ مقدار به صورت زیر است:

$$\begin{aligned} nT(X) &= \sum_{i=1}^n \|x_i - \hat{x}\|^2 \\ &= \sum_{i=1}^n \|x_i - \hat{x}\|^2 = \sum_{i=1}^n \|x_i - \mu_j + \mu_j - \hat{x}\|^2 \\ &= \sum_{i=1}^n \|x_i - \mu_j\|^2 + \sum_{i=1}^n \|\mu_j - \hat{x}\|^2 + \sum_{i=1}^n (x_i - \mu_j)^T (\mu_j - \hat{x}) \\ &= \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|^2 + \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|\mu_j - \hat{x}\|^2 + \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} (x_i - \mu_j)^T (\mu_j - \hat{x}) \end{aligned}$$

حال داریم :

$$\sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|^2 = \sum_{j=1}^k n_j W_j(X) \quad \text{where } n_j = \sum_{i=1}^n \gamma_{ij}$$

$$\sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|\mu_j - \hat{x}\|^2 = \sum_{j=1}^k n_j \|\mu_j - \hat{x}\|^2 = nB(X)$$

$$\mu_j = \frac{\sum_{i=1}^n \gamma_{ij} x_i}{\sum_{i=1}^n \gamma_{ij}} \implies \sum_{i=1}^n \gamma_{ij} (x_i - \mu_j)^T (\mu_j - \hat{x}) = 0$$

پس داریم :

$$nT(X) = \sum_{j=1}^k n_j W_j(X) + nB(X)$$

چون $T(X)$ ثابت است الگوریتم k-means $W(X)$ را کمینه و در نتیجه $B(X)$ را بیشینه می کند.
سوال ۴ حل.

ابتدا محاسبات مربوط به هر حالت را انجام می دهیم:

مقدار حالت S_2 :

$$V(S_2) = 1 + 0.9 \cdot V(S_2) \Rightarrow V(S_2) = 10$$

مقدار حالت S_4 :

$$V(S_4) = -5 + 0.9 \cdot \max(V(S_2), 0.2V(S_4) + 0.8V(S_2))$$

از آنجا که $V(S_2) = 10$:

$$V(S_4) = -5 + 0.9 \cdot 10 = 4$$

مقدار حالت S_3 :

اکشن اول:

$$V(S_3) = -1 + 0.9(0.2V(S_4) + 0.8V(S_2)) = -1 + 0.9(0.2 \cdot 4 + 0.8 \cdot 10) = -1 + 0.9(8.8) = 6.92$$

اکشن دوم:

$$V(S_3) = -1 + 0.9(0.2V(S_3) + 0.8V(S_4)) \Rightarrow V(S_3) = -1 + 0.18V(S_3) + 0.9 \cdot 4 = -1 + 0.18V(S_3) + 3.6$$

$$\Rightarrow V(S_3) = 0.18V(S_3) + 2.6 \Rightarrow V(S_3) < 3$$

پس بهترین اکشن، اولی است:

$$V(S_3) = 6.92$$

مقدار حالت S_1 :

اکشن اول:

$$V(S_1) = 0 + 0.9(0.4V(S_1) + 0.1V(S_2) + 0.5V(S_3)) = 0.9(0.4V(S_1) + 1 + 3.46) = 0.36V(S_1) + 4.014$$

حل این معادله:

$$V(S_1) = 0.36V(S_1) + 4.014 \Rightarrow 0.64V(S_1) = 4.014 \Rightarrow V(S_1) = \frac{4.014}{0.64} \approx 6.27$$

اکشن دوم:

$$V(S_1) = 0 + 0.9(0.2V(S_2) + 0.8V(S_3)) = 0.9(2 + 5/536) = 6/78$$

پس مقدار بهینه:

$$V(S_1) = 6/78$$

در صورتی که فرض کنیم پاداش در هنگام ورود به حالت مقصد (sj) دریافت می‌شود، مقادیر نهایی به شکل زیر خواهد بود:

$$V(S_1) = 7/536, \quad V(S_2) = 10, \quad V(S_3) = 8/8, \quad V(S_4) = 10$$

سوال ۵ حل.

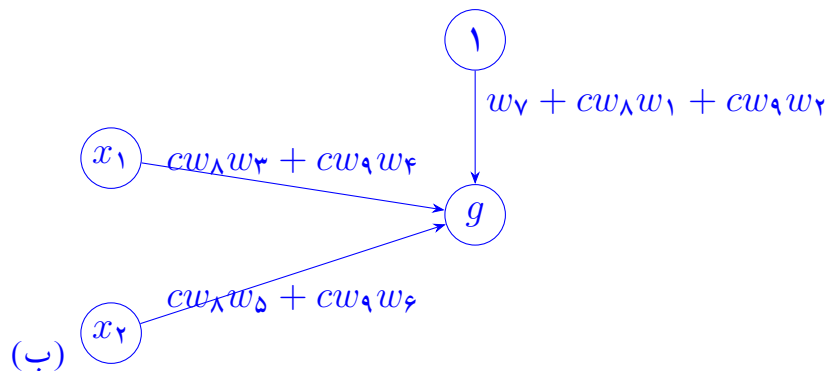
(I)

$$g(w_v + u_\lambda h(w_1 + w_3 x_1 + w_5 x_2) + u_9 h(w_2 + w_4 x_1 + w_6 x_2))$$

$$= \frac{1}{1 + \exp\left(-\left(w_v + cw_\lambda y_1 + cw_9 y_2 + (cw_\lambda w_3 + cw_9 w_4)x_1 + (cw_\lambda w_5 + cw_9 w_6)x_2\right)\right)}$$

The classification boundary is:

$$w_v + cw_\lambda w_1 + cw_9 w_2 + (cw_\lambda w_3 + cw_9 w_4)x_1 + (cw_\lambda w_5 + cw_9 w_6)x_2 = 0.$$



(ج) خیر چون شبکه جدید معادل شبکه اولیه است.

(د) Yes. If linear activation functions are used for all the hidden units, the output from hidden units can be written as a linear combination of the inputs; since these intermediate outputs serve as input for the final layer, there is an equivalent network without hidden layers.

سوال ۶ حل.

(الف)

برای اینکه نشان دهیم فرم مجموع مربعات وزنی با فرم ماتریسی معادل است، ابتدا باید خطای پیش‌بینی برای تمام نمونه‌ها را در یک بردار نمایش دهیم. این بردار خطا e از تفاضل بین مقادیر پیش‌بینی شده $(X\theta)$ و مقادیر واقعی (y) به دست می‌آید:

$$e = X\theta - y = \begin{pmatrix} \theta^T x^{(1)} - y^{(1)} \\ \theta^T x^{(2)} - y^{(2)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{pmatrix}$$

همانطور که مشخص است، هر درایه e_i از این بردار، همان خطای نمونه i -ام است. در قدم بعدی، برای اینکه وزن‌های $w^{(i)}$ را وارد محاسبات ماتریسی کنیم، یک ماتریس قطری W می‌سازیم که درایه‌های قطر اصلی آن همان وزن‌ها هستند و بقیه درایه‌ها صفرند.

$$W = \begin{pmatrix} w^{(1)} & 0 & \dots & 0 \\ 0 & w^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w^{(m)} \end{pmatrix}$$

حالا اگر عبارت ماتریسی $(X\theta - y)^T W (X\theta - y)$ را که معادل $e^T W e$ است بسط دهیم، به نتیجه زیر می‌رسیم:

$$\begin{aligned} e^T W e &= (e_1 \quad \dots \quad e_m) \begin{pmatrix} w^{(1)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & w^{(m)} \end{pmatrix} \begin{pmatrix} e_1 \\ \vdots \\ e_m \end{pmatrix} \\ &= (e_1 \quad \dots \quad e_m) \begin{pmatrix} w^{(1)} e_1 \\ \vdots \\ w^{(m)} e_m \end{pmatrix} \\ &= \sum_{i=1}^m w^{(i)} e_i^2 \end{aligned}$$

این عبارت نهایی دقیقاً همان فرم مجموع مربعات خطا است. بنابراین نشان دادیم که با تعریف ماتریس W به این شکل، دو عبارت با هم معادل هستند.

(ب)

در این قسمت، هدف ما پیدا کردن مقدار بهینه‌ی θ است که تابع هزینه $J(\theta)$ را کمینه کند. برای این کار، باید از تابع هزینه نسبت به θ گرادیان بگیریم و آن را برابر صفر قرار دهیم. تابع هزینه ما از قسمت قبل به این صورت است:

$$J(\theta) = \frac{1}{2} (X\theta - y)^T W (X\theta - y)$$

برای سادگی در مشتق‌گیری، ابتدا این عبارت را بسط می‌دهیم. پس از ضرب کردن جملات در هم، به عبارت زیر می‌رسیم:

$$J(\theta) = \frac{1}{2} (\theta^T X^T W X \theta - 2\theta^T X^T W y + y^T W y)$$

توجه کنید که دو جمله میانی در بسط اولیه با هم برابر بودند و در اینجا با هم ترکیب شدند. حالا از این عبارت گرادیان می‌گیریم. با استفاده از قواعد مشتق‌گیری ماتریسی، گرادیان هر جمله محاسبه می‌شود:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \frac{1}{2} [2(X^T W X)\theta - 2X^T W y + 0] \\ &= X^T W X \theta - X^T W y\end{aligned}$$

برای یافتن نقطه کمینه، گرادیان را برابر صفر می‌گذاریم و معادله را برای θ حل می‌کنیم:

$$X^T W X \theta = X^T W y$$

و در نهایت با فرض معکوس‌پذیر بودن ماتریس $X^T W X$ ، پاسخ نهایی که به معادله نرمال وزنی معروف است، به دست می‌آید:

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

این رابطه، فرم بسته برای یافتن پارامترهای بهینه در رگرسیون خطی وزنی است.

(ج)

در این بخش می‌خواهیم نشان دهیم که مسئله کمترین مربعات وزنی، از دیدگاه آماری با روش تخمین بیشینه درست‌نمایی (MLE) معادل است. فرض مسئله این است که هر $y^{(i)}$ از یک توزیع گوسی با میانگین $\theta^T x^{(i)}$ و واریانس $\sigma^{(i)2}$ می‌آید. برای شروع، تابع درست‌نمایی (Likelihood) کل داده‌ها را می‌نویسیم که حاصل ضرب احتمال هر یک از نمونه‌هاست:

$$L(\theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^{(i)2}}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^{(i)2}}\right)$$

کار کردن با ضرب سخت است، بنابراین از لگاریتم این تابع (Log-Likelihood) استفاده می‌کنیم که ضرب را به جمع تبدیل می‌کند و ماکزیمم کردن آن معادل ماکزیمم کردن خود تابع اصلی است.

$$\begin{aligned}\ell(\theta) = \log L(\theta) &= \sum_{i=1}^m \log \left[\frac{1}{\sqrt{2\pi\sigma^{(i)2}}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^{(i)2}}\right) \right] \\ &= -\frac{1}{2} \sum_{i=1}^m \log(2\pi\sigma^{(i)2}) - \frac{1}{2} \sum_{i=1}^m \frac{1}{\sigma^{(i)2}} (y^{(i)} - \theta^T x^{(i)})^2\end{aligned}$$

برای یافتن θ که این عبارت را ماکزیمم کند، متوجه می‌شویم که جمله اول (شامل لگاریتم) اصلاً به θ وابسته نیست و یک مقدار ثابت است. بنابراین، ماکزیمم کردن کل عبارت معادل کمینه کردن جمله دوم است:

$$\text{minimize}_{\theta} \left[\frac{1}{2} \sum_{i=1}^m \frac{1}{\sigma^{(i)2}} (y^{(i)} - \theta^T x^{(i)})^2 \right]$$

این عبارت دقیقاً مشابه تابع هزینه رگرسیون وزنی $J(\theta)$ است. با مقایسه این دو، به راحتی می‌توان دریافت که وزن هر نمونه $w^{(i)}$ با معکوس واریانس آن نمونه برابر است:

$$w^{(i)} = \frac{1}{\sigma^{(i)2}}$$

این نتیجه بسیار جالب است و نشان می‌دهد نمونه‌هایی که واریانس کمتری دارند (و در نتیجه قابل اعتمادتر هستند)، وزن بیشتری در مدل‌سازی خواهند داشت.