

تمرین تئوری ۱ یادگیری ماشین

نویسنده:

سید احسان حسن بیگی - ۴۰۲۲۱۱۷۲۳

بخش ۱

پرسش ۱

در الگوریتم یادگیری پرسپترون برای iteration t به ازای sample i از رخ داده داریم:

$$w^t = w^{t-1} + \gamma^{(i)} x^{(i)} \quad \text{update rule}$$

$$t \geq 1$$

بردار وزن را با تمام صفر initialize می‌کنیم

$$w^0 = 0$$

فرض داخلی داریم $w^t \cdot w^* = w^{t-1} \cdot w^* + \underbrace{\gamma^{(i)} x^{(i)T} w^*}_{\geq \gamma} \rightarrow w^t \cdot w^* \geq w^{t-1} \cdot w^* + \gamma$

طبق اثبات (I)

$$\left\{ \begin{array}{l} w^t \cdot w^* \geq t\gamma \\ \|w^t\| \|w^*\| \geq w^t \cdot w^* \end{array} \right. \rightarrow \textcircled{*} \quad \|w^t\| \geq t\gamma$$

بررسی norm 2 طبق update rule

$$\|w^t\|^2 = \|w^{t-1} + \gamma^{(i)} x^{(i)}\|^2 = \|w^{t-1}\|^2 + 2 \underbrace{\gamma^{(i)} x^{(i)T} w^{t-1}}_{\gamma^{(i)}} + \underbrace{\gamma^2 \|x^{(i)}\|^2}_{\leq R^2}$$

چون $\gamma^{(i)}$ و $\hat{\gamma}^{(i)}$ مختلف علامت هستند زیرا باید از تئوری ورنی شده است

$$\rightarrow \|w^t\|^2 \leq \|w^{t-1}\|^2 + R^2$$

طبق اثبات (II)

$$\textcircled{**} \quad \|w^t\|^2 \leq tR^2$$

از $\textcircled{*}$ و $\textcircled{**}$

$$t^2 \gamma^2 \leq \|w^t\|^2 \leq tR^2 \xrightarrow{t \geq 1} t \leq \left(\frac{R}{\gamma}\right)^2$$

اثبات I

استقرا؟ ثابت می کنیم که اگر $w^t \cdot w^* \geq t\gamma$ و $w^t \cdot w^* \geq w^{t-1} \cdot w^* + \gamma$ $t \geq 1$

نل 0 $t=1 \rightarrow w^1 \cdot w^* \geq \underbrace{w^0 \cdot w^*}_0 + \gamma \rightarrow w^1 \cdot w^* \geq \gamma \checkmark$

$t \rightarrow t+1$ $w^t \cdot w^* \geq t\gamma$
 $\left. \begin{array}{l} (y^{(i+k)} x^{(i+k)}) \cdot w^* \geq \gamma \end{array} \right\} + \rightarrow \underbrace{\left(w^t + y^{(i+k)} x^{(i+k)} \right)}_{w^{t+1}} \cdot w^* \geq (t+1)\gamma \checkmark$

اثبات II

استقرا؟ ثابت می کنیم که اگر $\|w^t\|^2 \leq tR^2$ و $\|w^t\|^2 \leq \|w^{t-1}\|^2 + R^2$ $t \geq 1$

نل 0 $t=1 \rightarrow \|w^1\|^2 \leq \underbrace{\|w^0\|^2}_0 + R^2 \rightarrow \|w^1\|^2 \leq R^2 \checkmark$

$t \rightarrow t+1$ $\|w^t\|^2 \leq tR^2$
 $\left. \begin{array}{l} \|y^{(i+k)} x^{(i+k)}\|^2 \leq R^2 \end{array} \right\} + \rightarrow \underbrace{\|w^t + y^{(i+k)} x^{(i+k)}\|^2}_{\|w^{t+1}\|^2} \leq \|w^t\|^2 + \|y^{(i+k)} x^{(i+k)}\|^2 \leq (t+1)R^2 \checkmark$

فرم ماتریسی

$$1 \quad J(w) = \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2 = \|Xw - y\|^2$$

$$X = \begin{bmatrix} -x^{(1)T} \\ \vdots \\ -x^{(n)T} \end{bmatrix}_{n \times (d+1)}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}_{n \times 1}$$

minimizing $J(w)$

$$\nabla_w J(w) = 2X^T(Xw - y) = 0$$

$$\rightarrow X^T X w = X^T y \rightarrow \hat{w} = \underbrace{(X^T X)^{-1} X^T}_{X^\dagger \text{ pseudo inverse}} y$$

2

- برای استفاده از این فرمول باید وارون ماتریس $X^T X$ را حساب کنیم که یک ماتریس $(d+1) \times (d+1)$ بعدی است و بسته به تعداد فیچرها می‌تواند هزینه‌بر باشد ← راه حل ← استفاده از روش‌های iterative مانند gradient descent
- از نظر تنوعی احتمال این وجود دارد که ماتریس $X^T X$ اصلاً وارون پذیر نباشد چون $\det(X^T X) = 0$ و البته اگر تعداد فیچرها از تعداد سمپل‌ها بیشتر باشد که قطعاً وارون پذیر نخواهد بود ← راه حل ← اگر از regularization استفاده کنیم عبارتی متناسب با λ در قطر اصلی ماتریس $X^T X$ اضافه می‌شود که کمک می‌کند $\det(X^T X) \neq 0$ و وارون پذیر شود
- ممکن است شرایط مسئله به صورتی باشد که تمام دنیا را در اول کار نداشته باشیم و stream سمپل‌ها در حین یادگیری دست ما برسد (online learning) ← راه حل ← در این حالت نیز استفاده از روش‌های iterative مانند gradient descent به کار می‌آید

$$\boxed{3} \quad J(w) = \sum_{i=1}^n F_i (y^{(i)} - w^T x^{(i)})^2 = (Xw - y)^T F (Xw - y)$$

$$\begin{bmatrix} x^{(1)} \cdot w - y^{(1)} & \dots & x^{(n)} \cdot w - y^{(n)} \end{bmatrix}_{1 \times n} \begin{bmatrix} F_1 (x^{(1)} \cdot w - y^{(1)}) \\ \vdots \\ F_n (x^{(n)} \cdot w - y^{(n)}) \end{bmatrix}_{n \times 1}$$

1) y, X

$$F = \begin{bmatrix} F_1 & 0 & \dots & 0 \\ 0 & F_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & F_n \end{bmatrix}_{n \times n}$$

$$\nabla_w J(w) = (Xw - y)^T \underbrace{(F + F^T)}_{2F} (X) = 0$$

$$\xrightarrow{T} X^T F^T (Xw - y) = 0 \rightarrow X^T F^T X w = X^T F^T y \rightarrow \hat{w} = (X^T F^T X)^{-1} X^T F^T y$$

$$= \boxed{(X^T F X)^{-1} X^T F y}$$

$$\boxed{4} \quad E_{n,y} [(y - w^T x)^2] = E_{n,y} [y^2 - 2w^T x y + (w^T x)^2]$$

$$= E_y [y^2] - 2E_{n,y} [w^T x y] + E_n [(w^T x)(x^T w)]$$

$$= E_y [y^2] - 2w^T \underbrace{E_{n,y} [x y]}_C + w^T \underbrace{E_n [x x^T]}_R w = E_y [y^2] - 2w^T C + w^T R w$$

$$\nabla_w E_{n,y} [(y - w^T x)^2] = R \hat{w} - C = 0 \rightarrow \text{نفسه } \hat{w} \text{ الة الة الة الة}$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y^{(i)} - \hat{y})^2$$

$\nwarrow \quad \nearrow$
 $w^T x^{(i)} \quad w^T x$

$$w^* = \underset{w}{\operatorname{argmin}} E_{n,y} [(y - \hat{y})^2] \xrightarrow{\frac{\partial}{\partial w}} \frac{\partial}{\partial w} E_{n,y} [(y - w^T x)^2]$$

$$= E_{n,y} [(y - w^* x) \cdot x] = 0 \quad (*)$$

فقطی که به واسطه تقریب وزن خط، بر اساس
دیتاست محدود داریم و اگر می‌دانستیم به ازای هر
X چه لایه‌ای باید داشته باشیم (دانستن $p(y|x)$)
خط بعدی می‌کشیدیم

\hat{w} وزن‌هایی است که با توجه به دیتاست محدود مان، آن را پیدا می‌کنیم
 w^* وزن‌هایی است که اگر توزیع $p(y|x)$ را داشته باشیم به آن خواهیم رسید
target distribution

$$\text{approximation error} = E_{n,y} [(y^* - \hat{y})^2] = E_{n,y} [(w^{*T} x - \hat{w}^T x)^2]$$

$$\text{structural error} = E_{n,y} [(y - y^*)^2] = E_{n,y} [(y - w^{*T} x)^2]$$

فقطی که به واسطه ساختار انتخاب شده (خط)
داریم و اگر ساختار بهینه‌ای پیدا می‌کردیم
می‌توانستیم در پیش‌بینی به حد دل‌برس
رسیدیم

$$\text{error} = E_{n,y} [(y - \hat{y})^2] = E_{n,y} [(y - \hat{w}^T x)^2] = E_{n,y} [(y - w^{*T} x + w^{*T} x - \hat{w}^T x)^2]$$

$$= \underbrace{E_{n,y} [(y - w^{*T} x)^2]}_{\text{structural error}} + \underbrace{E_{n,y} [(w^{*T} x - \hat{w}^T x)^2]}_{\text{approximation error}} + 2 E_{n,y} [(y - w^{*T} x)(w^{*T} x - \hat{w}^T x)]$$

$\underbrace{(w^{*T} - \hat{w}^T) x}_{x^T (w^* - \hat{w})}$ چون عدد است می‌توانیم T بگیریم

$$= \text{structural error} + \text{approximation error} + 2 E_{n,y} [(y - w^{*T} x) x^T (w^* - \hat{w})]$$

$\underbrace{\quad\quad\quad}_{0 \quad (*)}$
 $\quad\quad\quad 0$

1 $\hat{w} = (X^T X)^{-1} X^T y = \frac{1}{\underbrace{X^T X}_{n \times n}} \underbrace{X^T y}_{n \times 1} = \frac{\langle n, y \rangle}{\langle n, n \rangle}$ \otimes چون n بردار است داریم:

2 columns of X are orthogonal $\rightarrow X^T X = \begin{bmatrix} N & 0 & \dots & 0 \\ 0 & \langle n_1, n_1 \rangle & & \\ \vdots & & \ddots & \\ 0 & 0 & \dots & \langle n_d, n_d \rangle \end{bmatrix}$

$\rightarrow (X^T X)^{-1} X^T = \begin{bmatrix} \frac{1}{N} & & & \\ & \frac{1}{\langle n_1, n_1 \rangle} & & \\ & & \ddots & \\ & & & \frac{1}{\langle n_d, n_d \rangle} \end{bmatrix} \begin{bmatrix} 1 & 1 & \dots & 1 \\ n_1^{(1)} & n_1^{(2)} & \dots & n_1^{(N)} \\ \vdots & \vdots & \ddots & \vdots \\ n_d^{(1)} & n_d^{(2)} & \dots & n_d^{(N)} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} & \dots & \frac{1}{N} \\ \frac{n_1^{(1)}}{\langle n_1, n_1 \rangle} & \dots & \frac{n_1^{(N)}}{\langle n_1, n_1 \rangle} \\ \vdots & \ddots & \vdots \\ \frac{n_d^{(1)}}{\langle n_d, n_d \rangle} & \dots & \frac{n_d^{(N)}}{\langle n_d, n_d \rangle} \end{bmatrix}$

$\rightarrow (X^T X)^{-1} X^T y = \hat{w} = \begin{bmatrix} \frac{\sum_i y^{(i)}}{N} \\ \frac{\langle n_1, y \rangle}{\langle n_1, n_1 \rangle} \\ \vdots \\ \frac{\langle n_d, y \rangle}{\langle n_d, n_d \rangle} \end{bmatrix}$

\otimes مشاهده می‌کنیم که به برداری از فرایب
رگرسیون یک بعدی رسیدیم

3 $y = w_1 n + w_0 \rightarrow E[y] = E[w_1 n + w_0] = w_1 E[n] + w_0 \rightarrow w_0 = E[y] - w_1 E[n]$

$\text{cov}(n, y) = \text{cov}(n, w_1 n + w_0) = \text{cov}(n, w_1 n) + \text{cov}(n, w_0)$

$= w_1 \underbrace{\text{cov}(n, n)}_{\text{var}(n)} + w_0 \underbrace{\text{cov}(n, 1)}_0 \rightarrow \text{cov}(n, y) = w_1 \text{var}(n) \rightarrow w_1 = \frac{\text{cov}(n, y)}{\text{var}(n)}$

بخش ۲ (MLE, MAP)

پرسش ۱

$$D = \{x^{(1)}, \dots, x^{(N)}\}$$

likelihood

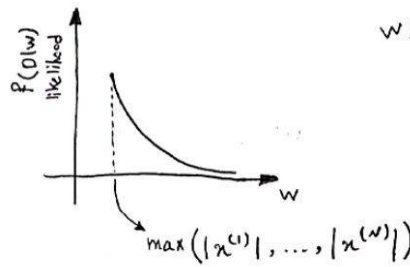
for continuous = $f(D|\theta)$
Random vars

$x^{(i)}$ are iid $U(-w, w) \rightarrow$

$$\begin{aligned} w &\geq \max(x^{(1)}, \dots, x^{(N)}) \\ -w &\leq \min(x^{(1)}, \dots, x^{(N)}) \end{aligned}$$

$$w \geq \max(|x^{(1)}|, \dots, |x^{(N)}|)$$

$$f(D|w) = \prod_{i=1}^N f_X(x^{(i)}|w) = \left(\frac{1}{2w}\right)^N$$



① مشاهده می شود که تابع likelihood نزولی است و اگر بخوانیم آنرا به سبب کنیم $\hat{w} = \max(|x^{(1)}|, \dots, |x^{(N)}|)$

1

$$\text{likelihood} = l(\mu) = p(D|\mu) = \prod_{i=1}^N f_X(x^{(i)}|\mu) = \prod_{i=1}^N \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x^{(i)}-\mu)^2}{2\sigma^2}}$$

$$\rightarrow L(\mu) = \ln l(\mu) = \sum_{i=1}^N -\ln(\sigma\sqrt{2\pi}) - \frac{(x^{(i)}-\mu)^2}{2\sigma^2}$$

$$\rightarrow \frac{\partial L(\mu)}{\partial \mu} = 0 \rightarrow \sum_{i=1}^N \frac{x^{(i)} - \mu}{\sigma^2} = 0 \rightarrow \sum_{i=1}^N x^{(i)} = N\mu \rightarrow \hat{\mu} = \frac{\sum_{i=1}^N x^{(i)}}{N}$$

2

فرض می‌کنیم $\text{prior} \sim N(\mu_0, \beta^2)$

$$p(D|\mu)p(\mu) = \prod_{i=1}^N f_X(x^{(i)}|\mu) \times p_\mu(\mu)$$

$$= \prod_{i=1}^N \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x^{(i)}-\mu)^2}{2\sigma^2}} \times \frac{1}{\beta\sqrt{2\pi}} e^{-\frac{(\mu-\mu_0)^2}{2\beta^2}}$$

$$\ln L(\mu) = \frac{(\mu-\mu_0)^2}{2\beta^2} - \ln(\beta\sqrt{2\pi})$$

$$\frac{\partial}{\partial \mu} = 0 \rightarrow \sum_{i=1}^N \frac{x^{(i)} - \mu}{\sigma^2} - \frac{1}{\beta^2} (\mu - \mu_0) = 0$$

$$\rightarrow \frac{N\mu}{\sigma^2} + \frac{\mu}{\beta^2} = \frac{1}{\sigma^2} \sum_{i=1}^N x^{(i)} + \frac{\mu_0}{\beta^2} \rightarrow (\mu) \left(\frac{N\beta^2 + \sigma^2}{\sigma^2\beta^2} \right) = \frac{\beta^2 \sum x^{(i)} + \sigma^2 \mu_0}{\sigma^2\beta^2}$$

$$\rightarrow \hat{\mu} = \frac{\beta^2 \sum x^{(i)} + \sigma^2 \mu_0}{N\beta^2 + \sigma^2} = \frac{\frac{\beta^2}{\sigma^2} \sum_{i=1}^N x^{(i)} + \mu_0}{\frac{\beta^2}{\sigma^2} N + 1}$$

3

② وقتی $N \rightarrow \infty$ در فرمول $\hat{\mu}_{MAP}$ برای μ_0 در صورت 1 در فرج بی اهمیت می‌شوند و به فرمول $\hat{\mu}_{MLE}$ می‌رسیم

یعنی وقتی تعداد نمونه‌ها بسیار زیاد باشند، هر دو تخمین یک جواب را می‌دهند و عملکرد می‌شوند

$$D = \{X^{(1)}, \dots, X^{(N)}\}$$

$$p(\beta | D) \propto \underbrace{p(D | \beta)}_{\text{iid } \prod_{i=1}^N p(x^{(i)} | \beta)} p(\beta)$$

$$p(D | \alpha, \beta) \rightarrow p(D | \beta) \quad \text{چون } \alpha \text{ معلوم است}$$

$$\begin{aligned} & \xrightarrow{\text{known}} \sim \text{gamma}(\alpha, \beta) \\ & \sim \text{gamma}(\alpha_0, \beta_0) \end{aligned}$$

$$p(D | \beta) = \prod_{i=1}^N \left(\frac{\beta^\alpha}{\Gamma(\alpha)} x^{(i)\alpha-1} e^{-\beta x^{(i)}} \right) = \frac{1}{\Gamma(\alpha)^N} \underbrace{\beta^{\alpha N} e^{-\beta \sum_{i=1}^N x^{(i)}}}_{(*)} \left(\prod_{i=1}^N x^{(i)} \right)^{\alpha-1} 1_{x^{(i)} > 0}$$

$$p(\beta) = \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \underbrace{\beta^{\alpha_0-1} e^{-\beta_0 \beta}}_{(**)} 1_{\beta > 0}$$

$$\xrightarrow{(*) (**)} p(\beta | D) \propto \beta^{\alpha N + \alpha_0 - 1} e^{-\beta(\beta_0 + \sum_{i=1}^N x^{(i)})} 1_{\beta > 0, x^{(i)} > 0}$$

$$\sim \text{gamma}(\alpha_0 + \alpha N, \beta_0 + \sum_{i=1}^N x^{(i)}) \quad \text{مناسب است}$$

⊗ بنابراین توزیع prior, posterior از یک جنس است، توزیع گاما برای توزیع گاما با α مشخص conjugate prior