

سوال ۱

a: فرض کنید یک داده به کلاس $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ تعلق داشته باشد و مقدار به دست آمده برای آن $\begin{bmatrix} 5 \\ 145 \\ 5 \end{bmatrix}$ باشد. مشخص است که مقدار این داده درست تعیین شده و به کلاس درستی تخصیص داده شده است. اما باز هم SSE بسیار بالایی خواهد داشت. بنابراین SSE اطلاعات خوبی برای $Classification$ نیست و تنها سعی می کند که داده ها را نزدیک خط جدا کننده نگه دارد. یعنی نه تنها به داده های نادرست و نویز خیلی حساس است، بلکه اگر یک داده به درستی در کلاس 1 باشد اما از صفحه جدا کننده فاصله زیادی داشته باشد، باز هم خطای خیلی زیادی SSE خروجی می دهد به صورتی که داد کاملاً دقیق و درست دسته بندی شده است. به عبارتی ممکن است تمام داده ها خیلی خوب و درست با الگوریتم دسته بندی شده باشند اما خطای SSE مقدار خیلی زیادی باشد.

b: برای رسیدن به جواب بهینه، از روش های $iterative$ مانند $gradient descent$ استفاده می شود. استفاده از $Perception Criterion$ یک فضای دیوید و محدب ($convex$) ایجاد می کند که می توان با استفاده از روش های $iterative$ به جواب بهینه رسید. اما اگر از $error-function$ به این صورت استفاده شود که تعداد داده های $misclassified$ را بشمارد، تقریباً همه جای فضای گمراهی منفرجه است و فضای صاف داریم. پس عملاً الگوریتم های $iterative$ نمی تواند به خط و خطی نود متوقف می شود چرا که گمراهی منفرجه است اما هنوز به جواب بهینه نرسیده است.

c: $Logistic regression$ یک روش احتمالی است در صورتی که $Perception$ یک تابع $Discriminant$ است. بنابراین با استفاده از $Logistic regression$ از مزایای مدل احتمالاتی بهره مند می شویم. مثلاً این که توانسته inference از $decision making$ جایی شود، می توانیم یک $rejection area$ تعریف کنیم و سایر مرزهای مدل احتمالاتی نسبت به مدل خیر احتمالاتی.

علاوه بر این ها، $Perception$ خیلی صفت به حالت چند کلاسه تقسیم می یابد در صورتی که در روش احتمالاتی تابع $softmax$ به راحتی جایگزین می شود. همچنین $Perception$ بین صفحه های جدا کننده که $margin$ معاف اما $error$ کمیانی داشته تعداد مائل نسبت به منجر به افزایش خطا

ردی داده های تست می شود. (به دلیل در نظر گرفتن margin)

اگر ملاحظه ما و هزینه misclassified کردن معض شود، در روش Perceptron سعی داریم با حذف اینها
توسط. اما در روش Logistic Regression، $P(L(x))$ به دست آمده همچنان معتبر است و باید فرایند
Decision making به ازای روش

d: در روش IRLS، محاسبه از ماتریس Hessian استفاده می شود. (با توجه به الگوریتم نیوتون-رافسون)

بنابراین هر step آن هزینه محاسباتی بهتری نسبت به step در GD دارد.
اما در هر مرحله، تخمین بهتری ارائه می دهد و در نهایت در تعداد مراحل کمتری جواب بهینه به دست می آید.
پس یک trade-off بین هزینه هر مرحله، تعداد مراحل ریجی که در برخی موارد GD در زمان کمتری
به جواب می رسد و در برخی موارد IRLS.

دلیل این که هر مرحله IRLS بهش از طولی کمتر این است که هر دفعه باید $(X^T X)^{-1}$ را حساب کند که
این کار هزینه بدتر است نسبت به محاسباتی گرایان. اما جابجایی در تعداد مراحل کمتری به جواب می رسد و
مکلف است در کل زمان رسیدن به جواب را کاهش دهد.

c: Probit Regression نسبت به outliers حساس تر است.

انتهای (tails) تابع Logistic Regression که همان Logistic Sigmoid است خیلی شبیه به
تابع $\exp(-x)$ عمل می کند. در صورتی که انتهای (tails) تابع Probit Regression خیلی شبیه به
 $x \rightarrow \infty$

تابع $\exp(-x^2)$ عمل می کند. بنابراین تابع Probit Regression احتمال خیلی کمی برای وجود outliers
 $x \rightarrow \infty$

قائل شده است و در صورت مشاهده آن، بسیار جریمه بالایی در نظر گرفته و در نتیجه به outliers حساس تر است.

۲ نقطه دلخواه باشد a ، b دارد J (یعنی متعلق به ناحیه J) در نظریه گریز. باید اثبات کنیم خط
 وصل کننده این ۲ نقطه نیز در ناحیه J قرار دارد.

$$\forall \lambda \in (0, 1) : x = \lambda a + (1 - \lambda)b$$

$$f_j(x) = f_j(\lambda a + (1 - \lambda)b) = \lambda f_j(a) + (1 - \lambda)f_j(b) \rightarrow \text{زیرا } f_j \text{ خطی است.}$$

$$\forall i \in \{1, 2, \dots, k\} : f_j(a) \geq f_i(a), f_j(b) \geq f_i(b) \rightarrow \text{زیرا } a, b \text{ در ناحیه } J \text{ است.}$$

$$\Rightarrow \lambda f_j(a) \geq \lambda f_i(a), (1 - \lambda)f_j(b) \geq (1 - \lambda)f_i(b)$$

$$\Rightarrow \forall i \in \{1, 2, \dots, k\} : \lambda f_j(a) + (1 - \lambda)f_j(b) \geq \lambda f_i(a) + (1 - \lambda)f_i(b)$$

$$\Rightarrow f_j(x) \geq f_i(x) \Rightarrow x \text{ نیز در ناحیه } J \text{ است.}$$

به ازای تمام نقاط x روی خط دامل a و b ، اثبات می شود که x هم «در ناحیه J است»
 پس decision regions شرط convex بودن را دارند. چون برای هر ناحیه دلخواه مانند J اثبات شد.

$$w^i = 0$$

سوال ۲

$$\begin{aligned} \forall x^j \rightarrow \text{misclassified}: w^{k+1} \cdot w^* &= (w^k + \epsilon x^j) \cdot w^* \\ &= w^k \cdot w^* + \epsilon (x^j \cdot w^*) \\ &= w^k \cdot w^* + \underbrace{w^{*T} x^j \epsilon}_{\geq \gamma} \geq w^k \cdot w^* + \gamma \end{aligned}$$

$$w^{k+1} \cdot w^* > w^k \cdot w^* + \gamma$$

$$w^k \cdot w^* > w^{k-1} \cdot w^* + \gamma$$

⋮

$$w^1 \cdot w^* > \gamma$$

به روشی مشابه داریم:

از استرا
دگرینی
نقته می شود

$$\boxed{w^{k+1} \cdot w^* > K\gamma} \rightarrow \text{پایان}$$

$$w^{k+1} \cdot w^* \leq \|w^{k+1}\| \|w^*\| = \|w^{k+1}\| \Rightarrow \boxed{\|w^{k+1}\| > K\gamma} \text{ lower bound}$$

$$\begin{aligned} \underbrace{w^{k+1} = w^k + \epsilon x^j}_{\text{update rule}} &\Rightarrow \|w^{k+1}\|^2 = \|w^k + \epsilon x^j\|^2 = \|w^k\|^2 + \epsilon^2 \|x^j\|^2 + 2\epsilon (w^k \cdot x^j) \\ &= \|w^k\|^2 + \epsilon^2 \|x^j\|^2 + 2\epsilon (w^k \cdot x^j) \epsilon^j \\ &\leq \|w^k\|^2 + \epsilon^2 \|x^j\|^2 \leq \|w^k\|^2 + R^2 \end{aligned}$$

$$\Rightarrow \|w^{k+1}\|^2 \leq \|w^k\|^2 + R^2 \xrightarrow[\text{فصل قبل}]{\text{از استرا مشابه}} \boxed{\|w^{k+1}\|^2 \leq K R^2} \text{ upper bound}$$

$$K\gamma^2 \leq \|w^{k+1}\|^2 \leq K R^2 \Rightarrow K\gamma^2 \leq K R^2 \Rightarrow \boxed{K \leq \frac{R^2}{\gamma^2}}$$

$$P(C_1) = \pi \Rightarrow P(C_r) = 1 - \pi$$

(سوال 4)

$$P(x|C_1) = N(x|\mu_1, \Sigma), P(x|C_r) = N(x|\mu_r, \Sigma), t_n \in \{1, 0\}$$

$$P(t|x, \mu_1, \mu_r, \Sigma, \pi) = \prod_{n=1}^N P(t_n|x_n) = \prod_{n=1}^N P(x_n|t_n) P(t_n)$$

$$= \prod_{n=1}^N (\pi N(x_n|\mu_1, \Sigma))^{t_n} ((1-\pi) N(x_n|\mu_r, \Sigma))^{1-t_n}$$

$$\Rightarrow -\ln P(t|x, \mu_1, \mu_r, \Sigma, \pi) = -\sum_{n=1}^N \{t_n \ln \pi + t_n \ln N(x_n|\mu_1, \Sigma) + (1-t_n) \ln(1-\pi) + (1-t_n) \ln N(x_n|\mu_r, \Sigma)\}$$

$$\Rightarrow \nabla_{\pi} \ln P(t|\mu_1, \mu_r, \Sigma, \pi) = \nabla_{\pi} \sum_{n=1}^N t_n \ln \pi + (1-t_n) \ln(1-\pi) = 0$$

$$\Rightarrow \sum_{n=1}^N \frac{t_n}{\pi} - \frac{1-t_n}{1-\pi} = \sum_{n=1}^N \frac{t_n - \ln \pi - \pi + t_n \pi}{\pi(1-\pi)} = \sum_{n=1}^N \frac{t_n - \pi}{\pi(1-\pi)} = 0$$

$$\Rightarrow \frac{1}{\pi(1-\pi)} \sum_{n=1}^N (t_n - \pi) = 0 \Rightarrow \sum_{n=1}^N t_n = N\pi \Rightarrow \boxed{\frac{1}{N} \sum_{n=1}^N t_n = \pi}$$

$$\nabla_{\mu_1} \ln P(t|\mu_1, \mu_r, \Sigma, \pi) = \nabla_{\mu_1} \sum_{n=1}^N t_n \ln N(x_n|\mu_1, \Sigma) = 0$$

$$\Rightarrow \nabla_{\mu_1} \sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) + \text{const} = 0 \Rightarrow \sum_{n=1}^N t_n x_n - \mu_1 = 0$$

$$\Rightarrow \sum_{n=1}^N t_n x_n = N\mu_1 \Rightarrow \boxed{\frac{1}{N} \sum_{n=1}^N t_n x_n = \mu_1} \xrightarrow{\text{class mean}} \boxed{\frac{1}{N} \sum_{n=1}^N (1-t_n) x_n = \mu_r}$$

ادامه سوال ۴)

$$\ln P(t | \mu_1, \mu_2, \Sigma, \pi) = \sum_{n=1}^N t_n \ln \pi + t_n \ln N(x_n | \mu_1, \Sigma) + (1-t_n) \ln(1-\pi) + (1-t_n) \ln N(x_n | \mu_2, \Sigma)$$

بجای های مربوط به Σ عبارت است از:

$$\begin{aligned} & \sum_{n=1}^N t_n \ln N(x_n | \mu_1, \Sigma) + (1-t_n) \ln N(x_n | \mu_2, \Sigma) \\ &= -\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) - \frac{1}{2} \sum_{n=1}^N (1-t_n) \ln |\Sigma| \\ & \quad - \frac{1}{2} \sum_{n=1}^N (1-t_n) (x_n - \mu_2)^T \Sigma^{-1} (x_n - \mu_2) + \text{Const} \end{aligned}$$

با حذف Const عبارت بالا داریم:

$$-\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr} \{ \Sigma^{-1} S \}, \quad S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n - \mu_1)(x_n - \mu_1)^T, \quad S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x_n - \mu_2)(x_n - \mu_2)^T$$

$$\Rightarrow \boxed{\sum^* = S}$$

$$P(C_1) = \frac{1}{N} \sum_{n=1}^N t_n, \quad P(C_2) = 1 - P(C_1) = \frac{1}{N} \sum_{n=1}^N (1-t_n)$$

خلاصه:

$$\mu_1 = \frac{1}{N} \sum_{n=1}^N t_n x_n, \quad \mu_2 = \frac{1}{N} \sum_{n=1}^N (1-t_n) x_n$$

$$\Sigma = \frac{1}{N} \left(\sum_{n \in C_1} (x_n - \mu_1)(x_n - \mu_1)^T + \sum_{n \in C_2} (x_n - \mu_2)(x_n - \mu_2)^T \right)$$

سؤال (۵)

داده ها به صورت خطی جدا پذیر هستند. پس داریم:

$$\forall x_i \in C_1: w^T x_i > 0$$

$$\forall x_i \in C_0: w^T x_i < 0$$

اگر مراد $a < 0$ را در w ضرب کنیم، روابط بالا همچنان برقرار خواهند بود. پس اگر w یک جواب درست باشد که داده ها را به درستی از هم جدا می کند، aw نیز یک جواب درست است.

بنابراین maximum likelihood داریم:

$$P(T|x) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{(1-t_n)}$$

$$\Rightarrow -\ln P(T|x) = -\sum_{n=1}^N t_n \ln y_n + (1-t_n) \ln (1-y_n)$$

$$\left\{ \begin{array}{l} t_n = 1 \Rightarrow w^T x_n > 0 \Rightarrow y_n = \frac{1}{1 + e^{-w^T x_n}} \rightarrow \text{یک عبارت مثبت است، پس اگر } a < 1 \text{ ضرب شود، } y_n \text{ نیز افزایش می یابد. پس } t_n \ln y_n \text{ نیز زیاد می شود.} \end{array} \right.$$

$$\left\{ \begin{array}{l} t_n = 0 \Rightarrow w^T x_n < 0 \Rightarrow y_n = \frac{1}{1 + e^{-w^T x_n}} \rightarrow \text{یک عبارت منفی است. پس اگر } a < 1 \text{ ضرب شود، } y_n \text{ کاهش می یابد. پس } (1-y_n) \ln (1-y_n) \text{ نیز زیاد می شود.} \end{array} \right.$$

بنابراین با ضرب w در $a < 1$ ، $\ln P(T|x)$ افزایش یافته پس $-\ln P(T|x)$ کاهش می یابد.

ادامه سؤال (۷)

بنابراین اگر w راه عددی بزرگتر از ضرب کنیم، یعنی فرض $1/w$ را از این دهیم، ادعا که همچنان داده ما را به دستی دسته بندی می کند، تا اینجا خطای ML آن کمتری شود. پس در حالت صدی، بهترین گزینه این است که w در صدها ضرب شود. بنابراین اندازه w به بی نهایت میل خواهد کرد.

که البته این اتفاق بدی است چرا که w بسیار بزرگ شده و محاسبه $verify$ می کند مسئله ای داده های $train$ و $test$ بدی شود.

سؤال ۶)

برای حل مدل رابطه صورت $(w^T x)$ در نظریه کرم چاکه $\text{logistic regression}$ هست.

پای داده های D_1 ، ویژگی های ما بردار x است پس به اندازه ی بعد x (که آن را d_1 می نامیم) پارامتر داریم. (یا مثلاً با در نظر گرفتن bias مدل پارامتر).

اما در مجموعه داده D_2 ، عمدتاً تعداد ویژگی ها برابر می شود. بنابراین تعداد پارامترهای w نیز برابر می شود. پس می توان گفت که زمان حل و میزان پردازش مورد نیاز برای دسته آردن پارامترهای همیشه نزدیکاً برابر می شود. (به عبارتی سرعت رسیدن به جواب نصف می شود.) از طرفی این ویژگی های جدید کاملاً مناسب هم هستند و اگر پارامترها کلی از مثلاً d_1 شروع کند، عمدتاً نصف پارامترها در هر مرحله برابر با نصف دیگر است. یعنی هیچ دیتای جدیدی به مسئله اضافه نشود و صرفاً هزینه بالارفته است.

سؤال (V)

$$\int_{-\infty}^{\infty} N(\theta|0,1) d\theta = 1, \quad N(\theta|0,1) \rightarrow \text{متوازن نسبت به 0}$$

$$\Rightarrow \int_{-\infty}^0 N(\theta|0,1) d\theta = \frac{1}{2}$$

$$\Phi(a) = \int_{-\infty}^a N(\theta|0,1) d\theta = \int_{-\infty}^0 N(\theta|0,1) d\theta + \int_0^a N(\theta|0,1) d\theta$$

$$= \frac{1}{2} + \int_0^a \frac{1}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}} d\theta \xrightarrow[\text{تغیر متغیر}]{\theta \rightarrow \sqrt{2}x} = \frac{1}{2} + \int_0^{\frac{a}{\sqrt{2}}} \frac{1}{\sqrt{\pi}} e^{-\frac{x^2}{2}} dx \sqrt{2}$$

$$= \frac{1}{2} + \int_0^{\frac{a}{\sqrt{2}}} \frac{1}{\sqrt{\pi}} e^{-\theta^2} d\theta = \frac{1}{2} + \frac{1}{2} \int_0^{\frac{a}{\sqrt{2}}} \frac{2}{\sqrt{\pi}} e^{-\theta^2} d\theta$$

$$= \frac{1}{2} \left(1 + \frac{2}{\sqrt{\pi}} \int_0^{\frac{a}{\sqrt{2}}} e^{-\theta^2} d\theta \right) = \frac{1}{2} \left(1 + \text{erf} \left(\frac{a}{\sqrt{2}} \right) \right)$$

$$\Rightarrow \Phi(a) = \frac{1}{2} \left(1 + \text{erf} \left(\frac{a}{\sqrt{2}} \right) \right)$$

$$\begin{cases} t_n = 1 & \text{if } d_n \geq \theta \\ t_n = 0 & \text{o.w.} \end{cases} \rightarrow f(d) = \int_{-\infty}^d G(\theta) d\theta, \text{ که } G(\theta) \text{ چگالی } \theta \text{ است.} \quad (\text{سؤال ۸})$$

$$P(t_n | x_n) = y_n^{t_n} (1 - y_n)^{1 - t_n}, \quad y_n = f(d_n) = \int_{-\infty}^{w^T x_n} G(\theta) d\theta$$

$$\Rightarrow P(T | X) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1 - t_n} \rightarrow -\ln P(T | X) = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln (1 - y_n)$$

$$\nabla_w E(w) = \nabla_y E(w) \times \nabla_w y = \sum_{n=1}^N \left\{ - \left(\frac{t_n}{y_n} - \frac{1 - t_n}{1 - y_n} \right) \nabla_w y_n \right\} \quad E(w)$$

$$= \sum_{n=1}^N \frac{y_n - t_n}{y_n (1 - y_n)} \nabla_w y_n = \sum_{n=1}^N \frac{(y_n - t_n) x_n}{y_n (1 - y_n)} \frac{df(w^T x_n)}{d(w^T x_n)} \rightarrow \text{رابطی کلاسیک}$$

$$\left\{ \begin{array}{l} t_i = 1 \rightarrow q_i = 1 \\ t_i = 0 \rightarrow q_i = -1 \end{array} \right\} \quad q_i = 2t_i - 1$$

متغیر q_i را اینجین کویف می‌کنیم:

$$\sum_{n=1}^N \frac{q_n x_n}{f(q_n w^T x_n)} \times \frac{df(q_n w^T x_n)}{d(q_n w^T x_n)}$$

پس رابطه به این شکل تبدیل می‌شود:

با توجه به این نکته $f(-a) = 1 - f(a)$ همچنین برای راضی $\frac{df(a)}{da}$ را به شکل $S(a)$ نمایش می‌دهیم.

پس برای Hessian داریم:

$$H = \nabla_w \nabla_w E(w) = \nabla_w \sum_{n=1}^N \frac{S(q_n w^T x_n) q_n}{f(q_n w^T x_n)} x_n$$

$$= \sum_{n=1}^N \nabla_w \left(\frac{S(q_n w^T x_n)}{f(q_n w^T x_n)} \right) q_n x_n = \sum_{n=1}^N \left\{ \frac{1}{f(q_n w^T x_n)} \times \nabla_w S(q_n w^T x_n) + S(q_n w^T x_n) \nabla_w \left(\frac{1}{f(q_n w^T x_n)} \right) \right\} q_n x_n$$

$$= \sum_{n=1}^N \left(\frac{-1}{f(q_n w^T x_n)} S(q_n w^T x_n) (q_n^T x_n) x_n^T - \frac{S(q_n w^T x_n)^2 q_n x_n^T}{f(q_n w^T x_n)^2} \right) q_n x_n$$

$$= \sum_{n=1}^N \left(\frac{S(q_n w^T x_n) (w^T x_n) q_n^T}{f(q_n w^T x_n)} + \frac{S(q_n w^T x_n)^2 q_n^T}{f(q_n w^T x_n)^2} \right) x_n x_n^T$$

$q_n \geq q_{n-1} : \text{sel}$

Hessian

اگر فرض کنیم که θ از (۱۰۵) به π می آید، خواهیم داشت:

$$\nabla_w L(w) = \sum_{n=1}^N \frac{(y_n - t_n)x_n}{y_n(1-y_n)} \times \frac{1}{\sqrt{2\pi}} e^{-\frac{(w^T x_n)^2}{2}}$$

$$\text{Hessian} = \sum_{n=1}^N \left\{ \frac{y_n^r \ln y_n - y_n \ln y_n}{y_n(1-y_n)} \times \frac{1}{\sqrt{r\pi}} e^{-\frac{r}{2}} - \frac{r \ln(y_n - t_n)}{y_n(1-y_n)} \times x_n x_n^T \times \frac{1}{\sqrt{r\pi}} e^{-\frac{r}{2}} \right\}$$

که دایم : ۱۷۷