

8.1

$$\begin{aligned}
 (a) \quad \frac{\partial b(a)}{\partial a} &= -\frac{1}{(1+e^{-a})^2} \cdot -e^{-a} \\
 &= \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{1+e^{-a}} \cdot \left(1 - \frac{1}{1+e^{-a}}\right) \\
 &= b(a) \cdot (1 - b(a))
 \end{aligned}$$

$$(b) \quad 1 - b(a) = \frac{e^{-a}}{1+e^{-a}} = \frac{1}{1+e^a} = b(-a)$$

$$\begin{aligned}
 (c) \quad y = b(a) &= \frac{1}{1+e^{-a}} \\
 y + ye^{-a} &= 1 \\
 e^{-a} &= \frac{1-y}{y} \\
 -a &= \log(1-y) - \log y \\
 a &= \log y - \log(1-y) \\
 b^{-1}(a) &= \log \frac{a}{1-a}
 \end{aligned}$$

8.2

$$\begin{aligned}
 (a) \quad E(w) &= \sum_i -\{y_i \log \pi_i + (1-y_i) \log (1-\pi_i)\} \\
 \nabla E(w) &= \sum_i -\left\{ \frac{y_i}{\pi_i} \cdot \pi_i (1-\pi_i) \cdot x_i - \frac{1-y_i}{1-\pi_i} \cdot \pi_i (1-\pi_i) \cdot x_i \right\} \\
 &= \sum_i -\{y_i (1-\pi_i) x_i - (1-y_i) \pi_i x_i\} \\
 &= \sum_i -\{y_i x_i - \pi_i x_i\} \\
 &= X \cdot (\pi - y)
 \end{aligned}$$

$$\begin{aligned}
 (b) \quad \nabla^2 E(w) &= \frac{\partial}{\partial w} \frac{\partial}{\partial w^T} [E(w)] \\
 &= \frac{\partial}{\partial w} (\pi - y)^T \cdot X^T \\
 &\quad \left(\text{consider } \frac{\partial}{\partial w} \pi^T \right) \\
 &= \left[\frac{\partial}{\partial w_i} \right]
 \end{aligned}$$

$$\begin{bmatrix} \vdots \\ \frac{\partial}{\partial w_n} \end{bmatrix} [\pi_1, \dots, \pi_n] = [\pi_1(1-\pi_1)x_1, \dots, \pi_n(1-\pi_n)x_n]$$

$$= X R$$

$$\nabla^2 E(w) = X R X^T$$

(c) since $\pi_i \in (0,1)$

Hence $R > 0$

$$(d) \quad w^{(new)} = w^{(old)} - [\nabla^2 E(w)]^{-1} \nabla E(w)$$

$$w^{(new)} = w^{(old)} - (X R X^T)^{-1} X (\pi - y)$$

$$= (X R X^T)^{-1} X R (X^T w^{(old)} - R^{-1}(\pi - y))$$

$$8.3 (a) \quad \nabla E(w) = 0$$

$$\Rightarrow (\pi - y) X = 0$$

$$\Rightarrow \pi - y = 0$$

$$\text{Hence when } y_i = 1 \Rightarrow \alpha(w^T x_i) \rightarrow 1$$

$$w^T x \rightarrow \infty$$

$$y_i = 0 \Rightarrow \alpha(w^T x_i) \rightarrow 0$$

$$w^T x \rightarrow -\infty$$

(b) If x is slightly across the margin

$$p(y|x) \rightarrow 0/1$$

It leaves no space for error/uncertainty.

8.4

$$\begin{aligned}
 (a) \quad \hat{E}(w) &= E(w) - \log p(w) \\
 &= E(w) - \left(\log \left(\frac{1}{(2\pi)^{\frac{1}{2}} |R|^{\frac{1}{2}}} \right) - \frac{1}{2} w^T P w \right) \\
 &= E(w) + \frac{1}{2} w^T P w
 \end{aligned}$$

$$\begin{aligned}
 (b) \quad \nabla \hat{E}(w) &= \nabla E(w) + P w \\
 &= X(\pi - y) + P w
 \end{aligned}$$

$$\begin{aligned}
 (c) \quad \nabla^2 \hat{E}(w) &= \nabla^2 E(w) + \nabla^2 \frac{1}{2} w^T P w \\
 &= X R X^T + \frac{\partial}{\partial w} \frac{\partial}{\partial w^T} \left(\frac{1}{2} w^T P w \right) \\
 &= X R X^T + \frac{\partial}{\partial w} w^T P^T \\
 &= X R X^T + P^T = X R X^T + P
 \end{aligned}$$

$$\begin{aligned}
 (d) \quad w^{(new)} &= w^{(old)} - (\nabla^2 \hat{E}(w))^{-1} \cdot \nabla \hat{E}(w) \\
 &= w^{(old)} - (X R X^T + P)^{-1} \cdot (X(\pi - y) + P w^{(old)}) \\
 &= (X R X^T + P)^{-1} (X R X^T + P) w^{(old)} - (X R X^T + P)^{-1} (X(\pi - y) + P w^{(old)}) \\
 &= (X R X^T + P)^{-1} (X R X^T w^{(old)} - X(\pi - y)) \\
 &= (X R X^T + P)^{-1} \cdot X R \cdot z
 \end{aligned}$$

$$\begin{aligned}
 (f) \quad \tilde{E}(w) &= E(w) + \frac{1}{2} w^T P w \\
 \text{if } P &= P_2, \text{ the last term, i.e. } \tilde{b} \\
 &\text{will not be constraint.}
 \end{aligned}$$

Consider (b) \Rightarrow gradient and analysis

8.5

(a) Define N_{mis} as the number of misclassification points

$$\text{By definition, } L_0(z_i) = \begin{cases} 0, & z_i \geq 0 \\ 1, & z_i < 0 \end{cases}$$

gives 0 when z_i is successfully predicted and

1 when z_i is misclassified, and $\sum_i L_0(z_i)$

is the total number of misclassified points

Hence $\min(N_{\text{mis}})$

$$= \min \left\{ \sum_i L_0(z_i) \right\}$$

(b) $R_{\text{emp}}(w)$

$$= \sum_{i \in u} -y_i w^T x_i$$

$$= \sum_{i \in u} -z_i$$

where u is set of misclassified points

Hence

$$L_p(z_i) = \begin{cases} 0, & z_i \geq 0 \\ -z_i, & z_i < 0 \end{cases}$$

$$= \max\{0, -z_i\}$$

(c) $R_{\text{emp}}(w)$

$$= \sum_i (y_i - w^T x_i)^2$$

$$= \sum_i (y_i^T y_i - 2y_i^T w^T x_i + x_i^T w w^T x_i)$$

$$= \sum_i (1 - 2z + x_i^T w y_i y_i^T w x_i)$$

$$= \sum_i (1 - 2z + z^T z)$$

$$= \sum_i (z_i - 1)^2$$

Here $L_{\text{log}}(z_i) = (z_i - 1)^2$

(4) $R_{\text{emp}}(w)$

$$= -\sum_i y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i)$$

$$= -\sum_{i \in I_1} y_i \log \pi_i + \sum_{i \in I_2} (1 - y_i) \log (1 - \pi_i)$$

where I_1 are those points with $y_i = 1$ (class label 1)

I_2 are those points with $y_i = 0$ (class label -1)

$R_{\text{emp}}(w)$

$$= -\sum_{i \in I_1} \log \pi_i + \sum_{i \in I_2} \log (1 - \pi_i)$$

$$= -\sum_{i \in I_1} \log \alpha(w^T x_i) + \sum_{i \in I_2} \log (1 - \alpha((1 - y_i) w^T x_i))$$

$$= -\sum_{i \in I_1} \log \alpha(z_i) + \sum_{i \in I_2} \log \alpha(y_i w^T x_i - w^T x_i)$$

$$= -\sum_{i \in I_1} \log \alpha(z_i) + \sum_{i \in I_2} \log \alpha(-w^T x_i)$$

$$= -\left(\sum_{i \in I_1} \log \alpha(z_i) + \sum_{i \in I_2} \log \alpha(z_i) \right)$$

$$= -\sum_i \log \alpha(z_i)$$

$$= \log (1 + e^{-z_i})$$

$$\propto \frac{1}{\log 2} \log (1 + e^{-z_i})$$

logistic regression



Intuitively,

- ① least-square tends to penalize the too correct output which is not desirable.
- ② 0-1 loss function is difficult to optimize which is not desirable.
- ③ perception and logistic regression are similar, however, logistic regression has some loss for correctly classified points near the boundary, which tends to push the boundary further away from those points.

Hence logistic regression may be better