

$$1. (1) \theta_{MAP} = \arg \max_{\theta} p(\theta | x, y) = \arg \max_{\theta} \frac{p(y|x, \theta) p(\theta | x)}{p(y|x)} = \arg \max_{\theta} p(y|x, \theta) p(\theta)$$

$$(2) \text{ 假设维数为 } n, \text{ 则 } p(\theta) = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2\sigma^2} \theta^T I_n \theta} = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2\sigma^2} \|\theta\|_2^2}$$

$$\Rightarrow \log p(\theta) = -\frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \|\theta\|_2^2 \propto -\frac{1}{2\sigma^2} \|\theta\|_2^2$$

此时 $\theta_{MAP} = \arg \max_{\theta} p(y|x, \theta) p(\theta) = \arg \max_{\theta} (\log p(y|x, \theta) + \log p(\theta))$

$$= \arg \max_{\theta} (\log p(y|x, \theta) - \frac{1}{2\sigma^2} \|\theta\|_2^2) = \arg \min_{\theta} (-\log p(y|x, \theta) + \frac{1}{2\sigma^2} \|\theta\|_2^2)$$

即为所求。比较可知此处 $\lambda = \frac{1}{2\sigma^2}$ 。

$$(3) \text{ 设 } X \in \mathbb{R}^{n \times n}, \theta_{MAP} = \arg \max_{\theta} p(\theta | X, \tilde{y}) = \arg \max_{\theta} \frac{p(\tilde{y} | X, \theta) p(\theta)}{p(\tilde{y} | X)} = \arg \max_{\theta} p(\tilde{y} | X, \theta) p(\theta)$$

由 $\varepsilon^{(i)} \sim N(0, \sigma^2)$, 有 $y^{(i)} \sim N(\theta^T x^{(i)}, \sigma^2)$, $p(\tilde{y}^{(i)} | x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2}$

$$\Rightarrow p(\tilde{y} | X, \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2} = \frac{1}{(\sqrt{2\pi\sigma^2})^n} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2} = \frac{1}{(\sqrt{2\pi\sigma^2})^n} e^{-\frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2}$$

由 $\theta \sim N(0, \eta^2 I_n)$ 有 $p(\theta) = \frac{1}{(\sqrt{2\pi\eta^2})^n} e^{-\frac{1}{2\eta^2} \|\theta\|_2^2}$

$$\Rightarrow \theta_{MAP} = \arg \max_{\theta} \frac{1}{(\sqrt{2\pi\sigma^2})^n} e^{-\frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2} \cdot \frac{1}{(\sqrt{2\pi\eta^2})^n} e^{-\frac{1}{2\eta^2} \|\theta\|_2^2}$$

$$= \arg \max_{\theta} \left(-\frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2 - \frac{n}{2} \log 2\pi - \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$= \arg \min_{\theta} \left(\frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$\text{令 } L(\theta) = \frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2, \nabla L(\theta) = \frac{1}{\sigma^2} X^T (X\theta - \tilde{y}) + \frac{1}{\eta^2} \theta$$

$$\text{令 } \nabla L(\theta) = 0, \text{ 有 } X^T X\theta + \frac{\eta^2}{\sigma^2} \theta - X^T \tilde{y} = 0 \Rightarrow (X^T X + \frac{\eta^2}{\sigma^2} I_n) \theta = X^T \tilde{y}$$

$$\Rightarrow \text{闭式解为 } \theta_{MAP} = (X^T X + \frac{\eta^2}{\sigma^2} I_n)^{-1} X^T \tilde{y}$$

$$(4) \theta_i \sim L(0, \lambda) \Rightarrow p(\theta_i) = \frac{1}{2\lambda} e^{-\frac{1}{2\lambda} |\theta_i|} \Rightarrow p(\theta) = \prod_{i=1}^n p(\theta_i) = \frac{1}{(2\lambda)^n} e^{-\frac{1}{2\lambda} \|\theta\|_1} = \frac{1}{(2\lambda)^n} e^{-\frac{1}{2\lambda} \|\theta\|_1}$$

与 (3) 相似, 可得到 $\theta_{MAP} = \arg \max_{\theta} \frac{1}{(\sqrt{2\pi\sigma^2})^n} e^{-\frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2} \cdot \frac{1}{(2\lambda)^n} e^{-\frac{1}{2\lambda} \|\theta\|_1}$

$$= \arg \min_{\theta} \left(\frac{1}{2\sigma^2} \|\tilde{y} - X\theta\|_2^2 + \frac{1}{2\lambda} \|\theta\|_1 \right)$$

$$= \arg \min_{\theta} (\|\tilde{y} - X\theta\|_2^2 + \lambda \|\theta\|_1)$$

即为所求。比较有 $\lambda = \frac{2\sigma^2}{\eta^2}$ 。

$$2. (1) \text{ 对正例, } \theta^T x = z > 0 \Rightarrow \lim_{c \rightarrow +\infty} p(y=1|x) = \lim_{c \rightarrow +\infty} \frac{1}{1 + e^{-c\theta^T x}} = \lim_{c \rightarrow +\infty} \frac{1}{1 + e^{-cE}} = 1$$

$$\text{对负例, } \theta^T x = z < 0 \Rightarrow \lim_{c \rightarrow +\infty} p(y=1|x) = \lim_{c \rightarrow +\infty} \frac{1}{1 + e^{-c\theta^T x}} = \lim_{c \rightarrow +\infty} \frac{1}{1 + e^{-cE}} = 0$$

$$(2) L(c, \theta) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} \log \frac{1}{1 + e^{-c\theta^T x}} + (1 - y^{(i)}) \log \frac{1}{1 + e^{c\theta^T x}})$$

$$= -\frac{1}{n} \sum_{i=1}^n (y^{(i)} \log (1 + e^{-c\theta^T x}) + (1 - y^{(i)}) \log (1 + e^{c\theta^T x}))$$

对每一个 i , 若 $y^{(i)} = 1 \Leftrightarrow \theta^T x > 0, \log (1 + e^{-c\theta^T x}) \rightarrow \log 1 = 0 \Rightarrow$ 求和由此项为 0

若 $y^{(i)} = 0 \Leftrightarrow \theta^T x < 0, \log (1 + e^{c\theta^T x}) \rightarrow \log 1 = 0 \Rightarrow$ 求和由此项为 0

$$\Rightarrow \lim_{c \rightarrow +\infty} L(c, \theta) = 0 \Rightarrow L(c, \theta) \text{ 随 } c \text{ 增大无限下降至 } 0.$$

(3) 由于数据不可分, 即使任意超平面均不能完全分开这两个类, 故无论如何调整 θ , 总有部分样本的预测无法达到其真实标签的极限, 故无法同时优化所有样本的预测概率。因此存在一个有限的最小损失值。由于有下界时必有下确界, 因此存在 $L_{min} \in \mathbb{R}$ s.t. $\forall \varepsilon, \exists \theta \in \mathbb{R}^d$ s.t. $L_{min} \leq L(\theta) \leq L_{min} + \varepsilon, \forall \theta \in \mathbb{R}^d, L_{min} \leq L(\theta)$ 。特别地, 取 $\varepsilon = \delta$, 记 θ_{δ} 符合条件, 则取 $\theta = \theta_{\delta}$ 并令 $n \rightarrow +\infty$, 即有损失函数能够收敛至 L_{min} 。

- (4) ① 不能。不论学习率高低, θ 均会发散至无穷远处, 无法收敛到收敛点。
 ② 能。由于 $\frac{1}{n} \sum_{i=1}^n \frac{1}{c} < +\infty$, 此时快速减小的学习率使得算法更稳定, 可能表现出更好的收敛性。
 ③ 不能。缩放不改变数据集的可分性, 仍然可能发散。
 ④ 能。 L^2 正则化可以限制 θ 的模长防止其无限增大, 使优化过程稳定并收敛至有限值。
 ⑤ 有可能。添加 Gauss 噪声可能导致数据集不再可分, 此时 $L(\theta)$ 有某一下界, 此时可能收敛。

$$3. (1) a_1 = \frac{e^{2i}}{e^{2i} + e^{2j} + e^{2k}} = 0.0900 \quad a_2 = 0.2447 \quad a_3 = 0.6653 \Rightarrow a^T = \begin{bmatrix} 0.0900 \\ 0.2447 \\ 0.6653 \end{bmatrix}$$

$$(2) NLL(0, y) = -y_1 \log a_3 = 1.6094$$

$$(3) NLL(a^T, y) = -\sum_{i=1}^n y_i \log a_i^T = -\sum_{i=1}^n y_i \log \left(\frac{e^{\frac{a_i^T}{\sum_{j=1}^n a_j^T}}}{\sum_{j=1}^n e^{\frac{a_j^T}{\sum_{j=1}^n a_j^T}}} \right) = -\sum_{i=1}^n y_i \left(\frac{a_i^T}{\sum_{j=1}^n a_j^T} - \log \sum_{j=1}^n e^{\frac{a_j^T}{\sum_{j=1}^n a_j^T}} \right)$$

由 y 是 one-hot 编码, 故 $\exists! j$ s.t. $y_j = 1$, 其余 $y_i = 0 \Leftrightarrow y_i = \delta_{ij}$
 故 $NLL(a^T, y) = -\frac{a_j^T}{\sum_{k=1}^n a_k^T} + \log \left(\sum_{k=1}^n e^{\frac{a_k^T}{\sum_{j=1}^n a_j^T}} \right)$ 。首先有 $\frac{\partial NLL(a^T, y)}{\partial a_j^T} = \frac{1}{\sum_{k=1}^n a_k^T} - \frac{\frac{a_j^T}{\sum_{k=1}^n a_k^T}}{\sum_{k=1}^n e^{\frac{a_k^T}{\sum_{j=1}^n a_j^T}}} = \frac{1}{\sum_{k=1}^n a_k^T} - \frac{a_j^T}{\sum_{k=1}^n a_k^T} = -\frac{a_j^T}{\sum_{k=1}^n a_k^T}$

由于 $a_j^T = \sum_{k=1}^n W_{kj}^T x_k + W_{0j}$, 有 $\frac{\partial a_j^T}{\partial W_{kj}} = \frac{\partial}{\partial W_{kj}} \sum_{k=1}^n W_{kj}^T x_k = x_k \delta_{jk}$

$$\frac{\partial NLL(a^T, y)}{\partial a_j^T} = -\delta_{jy} + \frac{e^{\frac{a_j^T}{\sum_{k=1}^n a_k^T}}}{\sum_{k=1}^n e^{\frac{a_k^T}{\sum_{j=1}^n a_j^T}}} = -\delta_{jy} + a_j^T$$

故 $\frac{\partial NLL(a^T, y)}{\partial a_j^T} = \frac{1}{\sum_{k=1}^n a_k^T} (-\delta_{jy} + a_j^T) x_k \delta_{jk} = (a_j^T - \delta_{jy}) x_k = x_k (a_j^T - y_j)$

因此 $\nabla_{W^T} NLL(a^T, y) = (x_k (a_j^T - y_j))_{k=1}^n = x (a^T - y)^T$

此处 $W^T = \begin{bmatrix} 1 & -1 & -2 \\ -1 & 2 & 1 \end{bmatrix}, W_0^T = [1, 1, 1], x = [1, 1, 1], z = (W^T x + W_0^T)^T = [0, 1, -1]^T$

$$a^T = SM(z^T) = [0.2447, 0.6653, 0.0900]^T, \quad y = [0, 1, 0]^T$$

有 $\nabla_{W^T} NLL(a^T, y) = x (a^T - y)^T = \begin{bmatrix} 0.2447 & -0.3348 & 0.0900 \\ 0.2447 & -0.3348 & 0.0900 \end{bmatrix}$

(4) 由 $a^T = [0.2447, 0.6653, 0.0900]^T$, 类别 1 对应第 2 个输出, 故概率为 66.52%

$$(5) W_{NEW}^T = W^T - \alpha \nabla_{W^T} NLL(a^T, y) = \begin{bmatrix} 1 & -1 & -2 \\ -1 & 2 & 1 \end{bmatrix} - 0.3 \begin{bmatrix} 0.2447 & -0.3348 & 0.0900 \\ 0.2447 & -0.3348 & 0.0900 \end{bmatrix}$$

$$= \begin{bmatrix} 0.8776 & -0.8326 & -2.0450 \\ -1.1224 & 2.1674 & 0.9550 \end{bmatrix}$$

$$(6) W_1^T = \begin{bmatrix} 0.8776 & -0.8326 & -2.0450 \\ -1.1224 & 2.1674 & 0.9550 \end{bmatrix} \Rightarrow z^T = (W_1^T x + W_0^T)^T = [-0.2447, 1.3348, -1.0900]^T$$

$$a_1^T = SM(z_1^T) = [0.1392, 0.7725, 0.0883]^T \Rightarrow \text{概率为 } 77.25\%$$

$$4. (1) \begin{bmatrix} z_1^T \\ z_2^T \\ z_3^T \\ z_4^T \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ 14 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ 13 \\ -4 \\ -15 \end{bmatrix}, \begin{bmatrix} a_1^T \\ a_2^T \\ a_3^T \\ a_4^T \end{bmatrix} = ReLU \begin{bmatrix} 2 \\ 13 \\ -4 \\ -15 \end{bmatrix} = \begin{bmatrix} 2 \\ 13 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} z_1^T \\ z_2^T \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & -1 & -1 & -1 \end{bmatrix} \begin{bmatrix} 2 \\ 13 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 13 \\ -13 \end{bmatrix}, \begin{bmatrix} a_1^T \\ a_2^T \end{bmatrix} = SM \begin{bmatrix} 13 \\ -13 \end{bmatrix} = \begin{bmatrix} 1.0000 \\ 0.6714 \times 10^{-13} \end{bmatrix}$$

习视为 $\begin{bmatrix} a_1^T \\ a_2^T \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ 为最终输出, $(f^T(z_1^T), f^T(z_2^T), f^T(z_3^T), f^T(z_4^T)) = (2, 13, 0, 0)$

$$(2) \text{ 对 } x^{(1)} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 1.2 \\ 1.3 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 1.0 \\ 1.1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \text{对 } x^{(1)} = \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 1.2 \\ 1.3 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1.0 \\ 1.1 \end{bmatrix} = \begin{bmatrix} -1 \\ -3 \end{bmatrix}$$

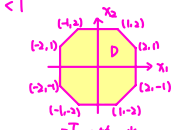
$$\begin{bmatrix} 1.0 \\ 1.1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow \text{对 } x^{(1)} = \begin{bmatrix} -3 \\ 0 \end{bmatrix}, \begin{bmatrix} 1.2 \\ 1.3 \end{bmatrix} = \begin{bmatrix} -4 \\ -0.3 \end{bmatrix}, \begin{bmatrix} 1.0 \\ 1.1 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \text{结果矩阵为 } \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix}$$

对 5 类为第二类, 即 $a_1^T < a_2^T$, 需要 $z_1^T < z_2^T \Leftrightarrow a_1^T + a_3^T + a_5^T + a_6^T < a_1^T - a_2^T - a_3^T - a_4^T + 2$
 $\Leftrightarrow a_1^T + a_1^T + a_3^T + a_6^T < 1 \Leftrightarrow \max \{z_1^T, 0\} + \max \{z_3^T, 0\} + \max \{z_5^T, 0\} + \max \{z_6^T, 0\} < 1$
 $\Leftrightarrow \max \{x_1 - 1, 0\} + \max \{x_2 - 1, 0\} + \max \{x_1 - 1, 0\} + \max \{x_2 - 1, 0\} < 1$

由于表达式较复杂, 故给出图像 (所有线均为直线)

对 $(x_1, x_2) \in D$, 分类结果为第二类

对 $(x_1, x_2) \in D^c$, 分类结果为第一类



(3) 对 Case 1, 由上述分析, 此处 $a_1^T + a_3^T + a_5^T + a_6^T = 0$, 输出为 $[0.1192, 0.8807]^T$ 为第二类

对 Case 2, $a_1^T + a_3^T + a_5^T + a_6^T = 1$, 输出为 $[0.5, 0.5]^T$ 无法分类

对 Case 3, $a_1^T + a_3^T + a_5^T + a_6^T = 3$, 输出为 $[0.9820, 0.0180]^T$ 为第一类。

$$(4) \frac{\partial \text{loss}}{\partial z_k^T} = a_k^T - y_k \Rightarrow \frac{\partial \text{loss}}{\partial z_1^T} = 1, \frac{\partial \text{loss}}{\partial z_2^T} = -1 \Rightarrow \frac{\partial \text{loss}}{\partial z^T} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \frac{\partial \text{loss}}{\partial W_0^T} = \frac{\partial \text{loss}}{\partial z^T} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\frac{\partial \text{loss}}{\partial W^T} = a^T \left(\frac{\partial \text{loss}}{\partial z^T} \right)^T = \begin{bmatrix} 2 \\ 13 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & -1 \end{bmatrix} = \begin{bmatrix} 2 & -2 \\ 13 & -13 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \frac{\partial \text{loss}}{\partial a^T} = W^T \frac{\partial \text{loss}}{\partial z^T} = \begin{bmatrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

$$\frac{\partial \text{loss}}{\partial z_j^T} = \frac{\partial \text{loss}}{\partial a_j^T} \cdot ReLU'(z_j^T) = \frac{\partial \text{loss}}{\partial a_j^T} \cdot \mathcal{I}_{(0,1)}(z_j^T) \Rightarrow \frac{\partial \text{loss}}{\partial z_1^T} = 2, \frac{\partial \text{loss}}{\partial z_2^T} = 2, \frac{\partial \text{loss}}{\partial z_3^T} = 0, \frac{\partial \text{loss}}{\partial z_4^T} = 0$$

$$\Rightarrow \frac{\partial \text{loss}}{\partial z^T} = \begin{bmatrix} 2 \\ 2 \\ 0 \\ 0 \end{bmatrix}, \frac{\partial \text{loss}}{\partial W_0^T} = \frac{\partial \text{loss}}{\partial z^T} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \frac{\partial \text{loss}}{\partial W^T} = x \left(\frac{\partial \text{loss}}{\partial z^T} \right)^T = \begin{bmatrix} 3 \\ 14 \end{bmatrix} \begin{bmatrix} 2 & 2 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 6 & 6 & 0 & 0 \\ 28 & 28 & 0 & 0 \end{bmatrix}$$

$$\Rightarrow \hat{W}^T = W^T - \eta \frac{\partial \text{loss}}{\partial W^T} = \begin{bmatrix} 0.8 & -0.8 \\ -0.3 & 0.3 \\ 1 & -1 \\ 1 & -1 \end{bmatrix}, \hat{W}_0^T = W_0^T - \eta \frac{\partial \text{loss}}{\partial W_0^T} = \begin{bmatrix} -0.1 \\ -1.2 \end{bmatrix}$$

$$\hat{W}_1^T = W_1^T - \eta \frac{\partial \text{loss}}{\partial W_1^T} = \begin{bmatrix} 0.4 & -0.6 & -1 & 0 \\ -2.8 & -1.8 & 0 & -1 \end{bmatrix}, \hat{W}_0^T = W_0^T - \eta \frac{\partial \text{loss}}{\partial W_0^T} = \begin{bmatrix} -1.2 \\ -1.2 \\ -1 \end{bmatrix}$$

以上即为更新后权重。