

第二次作业 胡延伸 PB22050983

1. 正则化的贝叶斯解释

(a) 证:

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta|x, y) = \underset{\theta}{\operatorname{argmax}} \frac{p(\theta, x, y)}{p(x, y)}$$

$$\text{其中, } \frac{p(\theta, x, y)}{p(x, y)} = \frac{p(\theta|y|x, \theta)}{p(x, y)} \cdot p(x, \theta)$$

由 $p(\theta) = p(\theta|x) \Rightarrow p(x, \theta) = p(x)p(\theta)$ 则:

$$\frac{p(\theta, x, y)}{p(x, y)} = \frac{p(y|x, \theta)}{p(x, y)} \cdot p(x)p(\theta) = \frac{p(y|x, \theta)p(\theta)}{p(y|x)}$$

由 $p(y|x)$ 与 θ 无关得: 当 $p(y|x, \theta)p(\theta)$ 最大, $p(y|x, \theta)p(\theta)/p(y|x)$ 最大

故 $\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|x, \theta)p(\theta)$

(b) 证:

$$\underset{\theta}{\operatorname{argmin}} -\log p(y|\theta, x) + \lambda \|\theta\|_2^2 = \underset{\theta}{\operatorname{argmax}} \frac{p(y|\theta, x)}{e^{\lambda \|\theta\|_2^2}}$$

由 (a) 得: $\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|\theta)p(\theta)$

$$\text{故只需证: } p(\theta) = e^{-\lambda \|\theta\|_2^2} = e^{-\lambda \theta^T \theta}$$

$$\text{由 } \theta \sim N(0, \eta^2 I) \Rightarrow p(\theta) = e^{-\frac{1}{2}\theta^T \theta / \eta^2}$$

$$\text{则 } \exists \lambda = \frac{1}{2\eta^2}, \text{ s.t. } \theta_{MAP} = \underset{\theta}{\operatorname{argmin}} -\log p(y|x, \theta) + \lambda \|\theta\|_2^2$$

(c) 证:

由题得: $\vec{y} = X^T \theta + E$, 其中 $E = (e^{(1)}, e^{(2)}, \dots)^T$, $e^{(i)}$ 相互之间独立同分布

$$\text{则 } E = y - X^T \theta \sim N(0, \sigma^2 I) \Rightarrow p(y|\theta, x) = e^{-\frac{1}{2\sigma^2} \|y - X^T \theta\|_2^2}$$

故由 (b) 得:

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmin}} -\log p(y|x, \theta) + \lambda \|\theta\|_2^2 = \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|y - X^T \theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2$$

(d) 证:

由题得:

$$X^T \theta - \vec{y} = E \sim N(0, \sigma^2 I) \text{ 即 } p(y|\theta, x) = e^{-\frac{1}{2\sigma^2} \|X^T \theta - \vec{y}\|_2^2}$$

又由 θ_i 独立同分布得: $p(\theta) = \prod p(\theta_i) = \left(\frac{1}{2\eta^2}\right)^n e^{-\frac{\|\theta\|_2^2}{\eta^2}}$, n 为 θ 的维度

由 (a) 得:

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|\theta, x)p(\theta) = \underset{\theta}{\operatorname{argmax}} e^{-\frac{1}{2\sigma^2} \|X^T \theta - \vec{y}\|_2^2 - \frac{1}{2\eta^2} \|\theta\|_2^2} = \underset{\theta}{\operatorname{argmin}} \left(\frac{1}{2\sigma^2} \|X^T \theta - \vec{y}\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$\Rightarrow \theta_{MAP} = \underset{\theta}{\operatorname{argmin}} \|X^T \theta - \bar{y}\|_2^2 + \frac{2\sigma^2}{b} \|\theta\|_1$$

即当 $\gamma = \frac{2\sigma^2}{b}$ 时, $\theta_{MAP} = \underset{\theta}{\operatorname{argmin}} J(\theta)$.

2. 逻辑回归: 训练稳定性

(a) 证:

$$\lim_{C \rightarrow \infty} P(y=1 | X^{(i)}) = \lim_{C \rightarrow +\infty} \frac{1}{1 + e^{-\theta^T X^{(i)}}}$$

对于正例, $\theta^T X^{(i)} > 0$ 则 $\lim_{C \rightarrow +\infty} \frac{1}{1 + e^{-\theta^T X^{(i)}}} = 1$;

对于正确分类的负例, $\theta^T X^{(i)} < 0$ 则 $\lim_{C \rightarrow +\infty} \frac{1}{1 + e^{-\theta^T X^{(i)}}} = 0$;

(b) 证:

① 对于正确分类的正例, $L(\theta) = -\log P(y=1 | X^{(i)}) \rightarrow 0 (C \rightarrow +\infty)$

对于正确分类的负例, $L(\theta) = -\log(1 - P(y=1 | X^{(i)})) \rightarrow 0 (C \rightarrow +\infty)$

则 $C \rightarrow +\infty$, $L(\theta) \rightarrow 0$ 即无限下降趋于 0.

② 在实际情况中, 这意味着 $\|\theta\|$ 不断增大, $L(\theta)$ 才能下降.

当参数更新步长太长, 会剧烈震荡; 当步长太小, 会导致收敛速度极慢.

(c)

当存在错误分类:

对于 $y^{(i)} = 0$, $\exists \theta^T X > 0$ 此时 $L(\theta) = -\log(1 - P(y=1 | X^{(i)})) \rightarrow +\infty (C \rightarrow +\infty)$

对于 $y^{(i)} = 1$, $\exists \theta^T X < 0$ 此时 $L(\theta) = -\log(P(y=1 | X^{(i)})) \rightarrow +\infty (C \rightarrow +\infty)$

$L(\theta)$ 显然存在大于 0 的下界, 因为错误分类的点对 $L(\theta)$ 的贡献恒 $\geq -\log 0.2$.

此时 $L(\theta)$ 取最优时, $\|\theta\|$ 为有限的值, 则 梯度下降会使得 θ 往最优点靠近.

(d)

(1) 应用不同学习率不会改变数据集线性可分这一性质, 故依然不会收敛

(2) 学习率下降会导致收敛速度更慢

(3) 线性缩放依然维持数据集线性可分, 故仍不会收敛.

(4) 可以收敛, 这样不会使得 $L(\theta)$ 最优时, $\|\theta\|$ 趋于无穷

(5) 可以收敛; 增加噪声可能使数据集变为线性不可分.

3. 多分类问题.

(a)

$$\vec{a} = SM(\vec{z}) = \left(\frac{e^{-1}}{e^{-1} + e^0 + e^1}, \frac{e^0}{e^{-1} + e^0 + e^1}, \frac{e^1}{e^{-1} + e^0 + e^1} \right) = (0.090, 0.245, 0.665)$$

(b)

$$NLL(\vec{a}, \vec{y}) = - \sum_{j=1}^{n_L} y_j \ln a_j^L = - (0, 0, 1) (ln 0.3, ln 0.5, ln 0.2)^T \approx 1.609$$

(c)

$$NLL = - \sum_{j=1}^{n_L} y_j \ln a_j^L = - \sum_{j=1}^{n_L} y_j \ln \frac{e^{z_j}}{\sum e^{z_i}} = - \sum_{j=1}^{n_L} y_j z_j + (\ln \sum e^{z_k}) \sum y_j$$

由于 \vec{y} 为 one-hot 编码则 $\sum y_j = 1$

$$\text{则: } NLL = - \sum_{j=1}^{n_L} y_j z_j + \ln \sum e^{z_k} \Rightarrow \frac{\partial NLL}{\partial W_{kj}} = \frac{\partial NLL}{\partial z_j} \cdot \frac{\partial z_j}{\partial W_{kj}} = (-y_j + a_j) x_k.$$

$$\text{对于, } W^L = \begin{bmatrix} 1 & -1 & -2 \\ -1 & 2 & 1 \end{bmatrix} \Rightarrow \vec{z} = xW = (0, 1, -1) \Rightarrow \vec{a} = (0.245, 0.665, 0.090)$$

$$\Rightarrow \nabla_{W^L} NLL = \begin{pmatrix} 0.245 & -0.335 & 0.090 \\ 0.245 & -0.335 & 0.090 \end{pmatrix}$$

(d)

$$\text{由 } \vec{a} = (0.245, 0.665, 0.090) \Rightarrow P_1 = 0.665$$

(e)

$$W^L = W^L - 0.5 \times \nabla_{W^L} NLL = \begin{pmatrix} 0.8775 & -1.3325 & -2.045 \\ -1.1225 & 1.6675 & 0.955 \end{pmatrix}$$

(f)

$$\vec{z} = xW = (-0.245, 0.335, -1.09) \Rightarrow \vec{a} = (0.311, 0.555, 0.134)$$

$$\Rightarrow P_1 = 0.555$$

4.

(a)

$$(z')^T = X^T W^1 + (W_0^1)^T = \begin{pmatrix} 3 & 14 \end{pmatrix} \begin{pmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{pmatrix} + (-1 & -1 & -1 & -1) = (2 \ 13 \ -4 \ -15)$$

$$\Rightarrow (a')^T = (2 \ 13 \ 0 \ 0) = (f'(z'_1), f'(z'_2), f'(z'_3), f'(z'_4))$$

$$\Rightarrow (z^2)^T = (a')^T W^2 + (W_0^2)^T = (2 \ 13 \ 0 \ 0) \begin{pmatrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{pmatrix} + (0 \ 2) = (15 \ -13)$$

$$\Rightarrow (a_1^2, a_2^2) = \left(\frac{e^{15}}{e^{15}+e^{-13}}, \frac{e^{-13}}{e^{15}+e^{-13}} \right) \approx (1.00, 0.00)$$

(b)

$$(z')^T = X^T W^1 + (W_0^1)^T = \begin{pmatrix} 0.5 & 0.5 \\ 0 & 2 \\ -3 & 0.5 \end{pmatrix} \begin{pmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{pmatrix} + (-1 & -1 & -1 & -1) = \begin{pmatrix} -0.5 & -0.5 & -1.5 & -1.5 \\ -1 & 1 & -1 & -3 \\ -4 & -0.5 & 2 & -1.5 \end{pmatrix}$$

$$\Rightarrow a'_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix}$$

(c)

(1)

$$(z^2)^T = (f'(z'_1), f'(z'_2), f'(z'_3), f'(z'_4)) \begin{pmatrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{pmatrix} + (0 \ 2) = (0 \ 2)$$

$$\Rightarrow (a_1^2, a_2^2) = \left(\frac{e^0}{e^0+e^2}, \frac{e^2}{e^0+e^2} \right)$$

(2)

$$\text{同上}, (z^2)^T = (1, -1) + (0, 2) = (1, 1) \Rightarrow (a_1^2, a_2^2) = (0.5, 0.5)$$

(3)

$$(z^2)^T = (3, -3) + (0, 2) = (3, -1) \Rightarrow (a_1^2, a_2^2) = \left(\frac{e^3}{e^3+e^{-1}}, \frac{e^{-1}}{e^3+e^{-1}} \right) \approx (0.982, 0.018)$$

(d)

$$\text{损失函数 } L = -y_1 \ln a_1^2 - y_2 \ln a_2^2 = -y_1 \ln \frac{e^{z_1^2}}{e^{z_1^2} + e^{z_2^2}} - y_2 \ln \frac{e^{z_2^2}}{e^{z_1^2} + e^{z_2^2}} = -y_1 z_1^2 - y_2 z_2^2 + (y_1 + y_2) \ln(e^{z_1^2} + e^{z_2^2})$$

$$\text{由 } y \text{ 为 one-hot 编码} \Rightarrow L = -y_1 z_1^2 - y_2 z_2^2 + \ln(e^{z_1^2} + e^{z_2^2})$$

由链式法则:

$$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial W^2} + \frac{\partial L}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial W^2} = \left(\frac{\partial L}{\partial z_1^2} \quad \frac{\partial L}{\partial z_2^2} \right) \left(\frac{\frac{\partial z_1^2}{\partial W^2}}{\frac{\partial z_2^2}{\partial W^2}} \right)$$

$$\text{由 } \frac{\partial L}{\partial z_1^2} = -z_1^2 + a_1^2 = -15 + 1 = 14, \quad \frac{\partial L}{\partial z_2^2} = -z_2^2 + a_2^2 = 13 + 0 = 13$$

$$\frac{\partial z_1^2}{\partial W^2} = \begin{pmatrix} a_1' & 0 \\ a_2' & 0 \\ a_3' & 0 \\ a_4' & 0 \end{pmatrix} \quad \frac{\partial z_2^2}{\partial W^2} = \begin{pmatrix} 0 & a_1' \\ 0 & a_2' \\ 0 & a_3' \\ 0 & a_4' \end{pmatrix}$$

$$\Rightarrow \frac{\partial L}{\partial W^2} = -14 \begin{pmatrix} 2 & 0 \\ 13 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} + 13 \begin{pmatrix} 0 & 2 \\ 0 & 13 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} -28 & 26 \\ -182 & 169 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}$$

$$\text{同理, } \frac{\partial L}{\partial W_0^2} = \left(\frac{\partial L}{\partial z_1^2}, \frac{\partial L}{\partial z_2^2} \right) \left(\frac{\partial z_1^2}{\partial W_0^2}, \frac{\partial z_2^2}{\partial W_0^2} \right) = -14 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 13 \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -14 \\ 13 \end{pmatrix}$$

$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial z^2} \cdot \frac{\partial(z_1^2, z_2^2)}{\partial(a_1', \dots, a_4')} \cdot \begin{pmatrix} \frac{\partial a_1'}{\partial z_1'} \\ \vdots \\ \frac{\partial a_4'}{\partial z_4'} \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial z_1'}{\partial W^1} \\ \vdots \\ \frac{\partial z_4'}{\partial W^1} \end{pmatrix}, \text{ 其中最后两项为元素相乘.}$$

$$\frac{\partial L}{\partial W_0^1} = \frac{\partial L}{\partial z^2} \cdot \frac{\partial(z_1^2, z_2^2)}{\partial(a_1', \dots, a_4')} \cdot \begin{pmatrix} \frac{\partial a_1'}{\partial z_1'} \\ \vdots \\ \frac{\partial a_4'}{\partial z_4'} \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial z_1'}{\partial W_0^1} \\ \vdots \\ \frac{\partial z_4'}{\partial W_0^1} \end{pmatrix}$$

$$\text{由 } \frac{\partial(z_1^2, z_2^2)}{\partial(a_1', \dots, a_4')} = (W^2)^T, \quad \begin{pmatrix} \frac{\partial a_1'}{\partial z_1'} \\ \vdots \\ \frac{\partial a_4'}{\partial z_4'} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix},$$

$$\frac{\partial z_1'}{\partial W^1} = \begin{pmatrix} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \end{pmatrix}, \quad \frac{\partial z_2'}{\partial W^1} = \begin{pmatrix} 0 & x_1 & 0 & 0 \\ 0 & x_2 & 0 & 0 \end{pmatrix}, \dots, \quad \frac{\partial z_4'}{\partial W^1} = \begin{pmatrix} 0 & 0 & 0 & x_1 \\ 0 & 0 & 0 & x_2 \end{pmatrix}$$

$$\frac{\partial z_1'}{\partial W_0^1} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad \dots, \quad \frac{\partial z_4'}{\partial W_0^1} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

得:

$$\frac{\partial L}{\partial W^1} = (-14 \ 13) \cdot \begin{pmatrix} 1 & 1 & 1 & 1 \\ -1 & -1 & -1 & -1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} \frac{\partial z_1'}{\partial W_0^1} \\ \vdots \\ \frac{\partial z_4'}{\partial W_0^1} \end{pmatrix} = -27 \begin{pmatrix} x_1 & x_2 & 0 & 0 \\ x_2 & x_2 & 0 & 0 \end{pmatrix} = \begin{pmatrix} -81 & -81 & 0 & 0 \\ -378 & -378 & 0 & 0 \end{pmatrix}$$

$$\frac{\partial L}{\partial W_0^1} = -27 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} -27 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\text{则更新后: } W^2 = \begin{pmatrix} 3.8 & -3.6 \\ 19.2 & -17.9 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad W_0^2 = \begin{pmatrix} 1.4 \\ 0.7 \end{pmatrix}$$

$$W^1 = \begin{pmatrix} 9.2 & 8.1 & 1 & 0 \\ 37.8 & 38.8 & 0 & -1 \end{pmatrix}, \quad W_0^1 = \begin{pmatrix} 1.7 \\ 1.7 \\ -1 \\ -1 \end{pmatrix}$$