

IML 第三次作业

习题 5.1

如果用线性函数，则每一层的输出都是上一层的线性模型，最终的输出也只是复杂的线性模型，无法表示复杂的非线性关系。

作业 5.2

当 x 很大时， $\exp(x)$ 的结果可能发生溢出而显示 NaN；假设 $x_j (1 \leq j \leq C)$ 的最大值为 x^* ，可以如下处理

$$\frac{\exp(x_i)}{\sum_{j=1}^C \exp(x_j)} = \frac{\exp(x_i) \exp(-x^*)}{\exp(-x^*) \sum_{j=1}^C \exp(x_j)} = \frac{\exp(x_i - x^*)}{\sum_{j=1}^C \exp(x_j - x^*)}$$

这样可以避免数值上溢问题。

同样地，对于第二个式子，可以如下处理

$$\log \sum_{j=1}^C \exp(x_j) = \log \left[\exp(x^*) \sum_{j=1}^C \exp(x_j - x^*) \right] = x^* + \log \sum_{j=1}^C \exp(x_j - x^*)$$

作业 5.3

$$\text{令 } f(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^C \exp(x_j)}, \quad g(x_i) = \log f(x_i)。$$

首先

$$\frac{\partial f(x_i)}{\partial \mathbf{x}} = \left[\frac{\partial f(x_i)}{\partial x_1} \quad \frac{\partial f(x_i)}{\partial x_2} \quad \cdots \quad \frac{\partial f(x_i)}{\partial x_C} \right]$$

当 $i = k$ 时，有

$$\frac{\partial f(x_i)}{\partial x_k} = \frac{\exp(x_k) \sum_{j \neq k}^C \exp(x_j)}{\left[\sum_{j=1}^C \exp(x_j) \right]^2}$$

当 $i \neq k$ 时，有

$$\frac{\partial f(x_i)}{\partial x_k} = - \frac{\exp(x_i) \exp(x_k)}{\left[\sum_{j=1}^C \exp(x_j) \right]^2}$$

其次

$$\frac{\partial g(x_i)}{\partial \mathbf{x}} = \left[\frac{\partial g(x_i)}{\partial x_1} \quad \frac{\partial g(x_i)}{\partial x_2} \quad \cdots \quad \frac{\partial g(x_i)}{\partial x_C} \right]$$

由链式法则（视 \log 以 e 为底）

$$\frac{\partial g(x_i)}{\partial x_k} = \frac{1}{f(x_i)} \cdot \frac{\partial f(x_i)}{\partial x_k}$$

当 $i = k$ 时, 有

$$\frac{\partial g(x_i)}{\partial x_k} = \frac{\sum_{j \neq k}^C \exp(x_j)}{\sum_{j=1}^C \exp(x_j)}$$

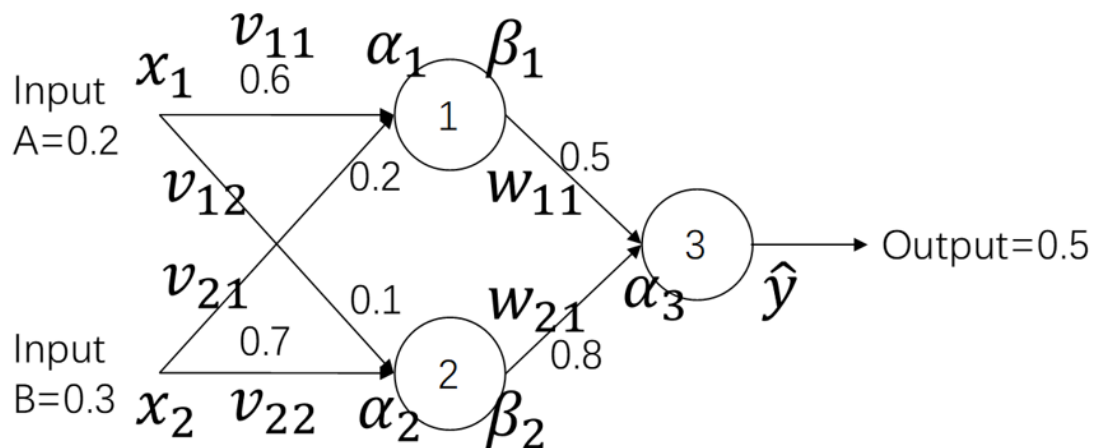
当 $i \neq k$ 时, 有

$$\frac{\partial g(x_i)}{\partial x_k} = -\frac{\exp(x_i)}{\sum_{j=1}^C \exp(x_j)}$$

作业 5.4

$\text{ReLU}(x) = \max(0, x)$, 其导数为

$$\text{ReLU}'(x) = \mathbb{I}(x > 0)$$



参数更新前, 各部分输入输出结果为

$$\alpha_1 = v_{11}x_1 + v_{21}x_2 = 0.18$$

$$\beta_1 = \text{ReLU}(\alpha_1) = 0.18$$

$$\alpha_2 = v_{12}x_1 + v_{22}x_2 = 0.23$$

$$\beta_2 = \text{ReLU}(\alpha_2) = 0.23$$

$$\alpha_3 = w_{11}\beta_1 + w_{21}\beta_2 = 0.274$$

$$\hat{y} = \text{ReLU}(\alpha_3) = 0.274$$

$$E = \frac{1}{2}(y - \hat{y})^2 = 0.026$$

计算梯度项

$$\frac{\partial E}{\partial v_{11}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) w_{11} \text{ReLU}'(\alpha_1) x_1 = -0.023$$

$$\frac{\partial E}{\partial v_{12}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) w_{11} \text{ReLU}'(\alpha_2) x_1 = -0.023$$

$$\frac{\partial E}{\partial v_{21}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) w_{21} \text{ReLU}'(\alpha_1) x_2 = -0.054$$

$$\frac{\partial E}{\partial v_{22}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) w_{21} \text{ReLU}'(\alpha_2) x_2 = -0.054$$

$$\frac{\partial E}{\partial w_{11}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) \beta_1 = -0.041$$

$$\frac{\partial E}{\partial w_{21}} = -(y - \hat{y}) \text{ReLU}'(\alpha_3) \beta_2 = -0.052$$

因为学习率 $\eta = 1$, 更新参数:

$$v_{11} \leftarrow v_{11} - \frac{\partial E}{\partial v_{11}} = 0.623$$

$$v_{12} \leftarrow v_{12} - \frac{\partial E}{\partial v_{12}} = 0.123$$

$$v_{21} \leftarrow v_{21} - \frac{\partial E}{\partial v_{21}} = 0.254$$

$$v_{22} \leftarrow v_{22} - \frac{\partial E}{\partial v_{22}} = 0.754$$

$$w_{11} \leftarrow w_{11} - \frac{\partial E}{\partial w_{11}} = 0.541$$

$$w_{21} \leftarrow w_{21} - \frac{\partial E}{\partial w_{21}} = 0.852$$

更新参数后

$$\alpha_1 = v_{11}x_1 + v_{21}x_2 = 0.201$$

$$\beta_1 = \text{ReLU}(\alpha_1) = 0.201$$

$$\alpha_2 = v_{12}x_1 + v_{22}x_2 = 0.251$$

$$\beta_2 = \text{ReLU}(\alpha_2) = 0.251$$

$$\alpha_3 = w_{11}\beta_1 + w_{21}\beta_2 = 0.221$$

$$\hat{y} = \text{ReLU}(\alpha_3) = 0.323$$

$$E = \frac{1}{2}(y - \hat{y})^2 = 0.016$$

可见参数更新使平方损失值下降了。

习题 6.4

线性判别分析可以解决多分类问题，而线性核支持向量机只能解决二分类问题。并且线性判别分析只有处理线性可分样本时才工作地比较好。所以二者在处理二分类问题时且两类样本线性可分时等价。

习题 6.6

SVM 训练出来的模型只和稀疏的支持向量有关，但是若噪声出现在这些支持向量中，就会参与到优化问题的最大化中，即对 SVM 的最终模型造成很大的影响。

习题 6.9

书中式 (6.29) 如下：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{0/1} \left(y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1 \right)$$

用对率损失 $\ell_{\log}(z) = \log(1 + \exp(-z))$ 替代 $\ell_{0/1}$ ，令 $z = y_i(\mathbf{w}^\top \mathbf{x}_i + b)$ ，则式 (6.29) 变为

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \log(1 + \exp(-z))$$

这是一个带 L_2 正则项的正则化问题，根据表示定理，其解总可以写成

$$\mathbf{w}_* = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i)$$

按书上对对偶问题的推导，其对偶问题为

$$\min_{\alpha} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) + C \sum_{i=1}^m \log \left[1 + \exp \left(-y_i \sum_{j=1}^m \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right) \right]$$

上面这个问题没有约束，可用 GD 算法求解后得到 \mathbf{w}_*, b_* 。

作业 6.4

令 $\alpha^* = \begin{pmatrix} \alpha \\ \hat{\alpha} \end{pmatrix}$ ，则

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i) (\alpha_j - \hat{\alpha}_j) \kappa(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j k_{ij} - \hat{\alpha}_i \alpha_j k_{ij} - \alpha_i \hat{\alpha}_j k_{ij} + \hat{\alpha}_i \hat{\alpha}_j k_{ij} \\ &= \alpha^{*\top} \begin{bmatrix} k & -k \\ -k & k \end{bmatrix} \alpha^* \end{aligned}$$

令 $\mathbf{v} = \begin{pmatrix} -\mathbf{y} - \epsilon \\ \mathbf{y} - \epsilon \end{pmatrix}$, 则有

$$\sum_{i=1}^m (y_i (\hat{\alpha}_i - \alpha_i) - \epsilon (\hat{\alpha}_i + \alpha_i)) = \boldsymbol{\alpha}^{*\top} \mathbf{v}$$

因此原式可化为

$$\max_{\boldsymbol{\alpha}^*} g(\boldsymbol{\alpha}^*) = \boldsymbol{\alpha}^{*\top} \mathbf{v} - \frac{1}{2} \boldsymbol{\alpha}^{*\top} \mathbf{K} \boldsymbol{\alpha}^*$$

$$\text{s.t. } C \succ \boldsymbol{\alpha}^* \succ 0, \boldsymbol{\alpha}^{*\top} \mathbf{v} = 0$$

这里 $\mathbf{K} = \begin{bmatrix} k & -k \\ -k & k \end{bmatrix}$ 。

作业 6.5

简单起见, 不妨令两个变量为 \mathbf{x}, \mathbf{y} , 并且他们都是 n 维的

$$\begin{aligned} \kappa(\mathbf{x}, \mathbf{y}) &= \left(\sum_{i=1}^n x_i y_i \right)^2 \\ &= \sum_{i=1}^n (x_i^2) (y_i^2) + \sum_{i=2}^n \sum_{j=1}^{i-1} (\sqrt{2} x_i x_j) (\sqrt{2} y_i y_j) \end{aligned}$$

所以

$$\phi(\mathbf{x}) = (x_n^2, \dots, x_1^2, \sqrt{2} x_n x_{n-1}, \dots, \sqrt{2} x_n x_1, \sqrt{2} x_{n-1} x_{n-2}, \dots, \sqrt{2} x_2 x_1)$$