

HWS.

5.1 $f(x) = W^T x$

线性激活函数在定义域内变换情况相同, 而理想中的激活函数是能够把可能在较大范围内变化的输入值挤入到 (0,1) 输出值范围内的函数。

2. 讨论 $\frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$ 和 $\log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$ 的数值溢出问题

实数在计算机中以二进制定表示, 数值过小会取 0 (下溢出), 数值过大会上溢出。令 $f(x) = \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$, $g(x) = \log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$

当 $x_i \rightarrow -\infty$, $i=1, 2, \dots, c$ 时 $\sum_{j=1}^c \exp(x_j) \rightarrow 0$, 易发生下溢出

当 $x_i \rightarrow +\infty$, $i=1, 2, \dots, c$ 时, $\exp(x_i)$ 可能发生上溢出

对 $f(x)$, 令 $X = \max(x_i)$, $i=1, 2, \dots, c$, 把 $f(x)$ 改为计算 $f(x_i - X)$

即可避免溢出

$$\frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)} = \frac{\exp(x_i - X)}{\sum_{j=1}^c \exp(x_j - X)} = \frac{e^{x_i - X}}{\sum_{j=1}^c e^{x_j - X}} \quad \left\{ \begin{array}{l} e^{x_i - X} \leq e^{M-M} = 1, \text{ 分子不会上溢出} \\ \sum_{j=1}^c e^{x_j - X} \geq e^{M-M} = 1, \text{ 分母不会下溢出} \end{array} \right.$$

$$\log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)} = \log \frac{\exp(x_i - X)}{\sum_{j=1}^c \exp(x_j - X)} = \log \frac{e^{x_i - X}}{\sum_{j=1}^c e^{x_j - X}} = x_i - X - \log \sum_{j=1}^c e^{x_j - X}$$

$$-1 = e^{x_i - X} \leq \sum_{j=1}^c e^{x_j - X} \leq \sum_{j=1}^c e^{X-X} = c$$

分子不会上溢出, 分母不会下溢出

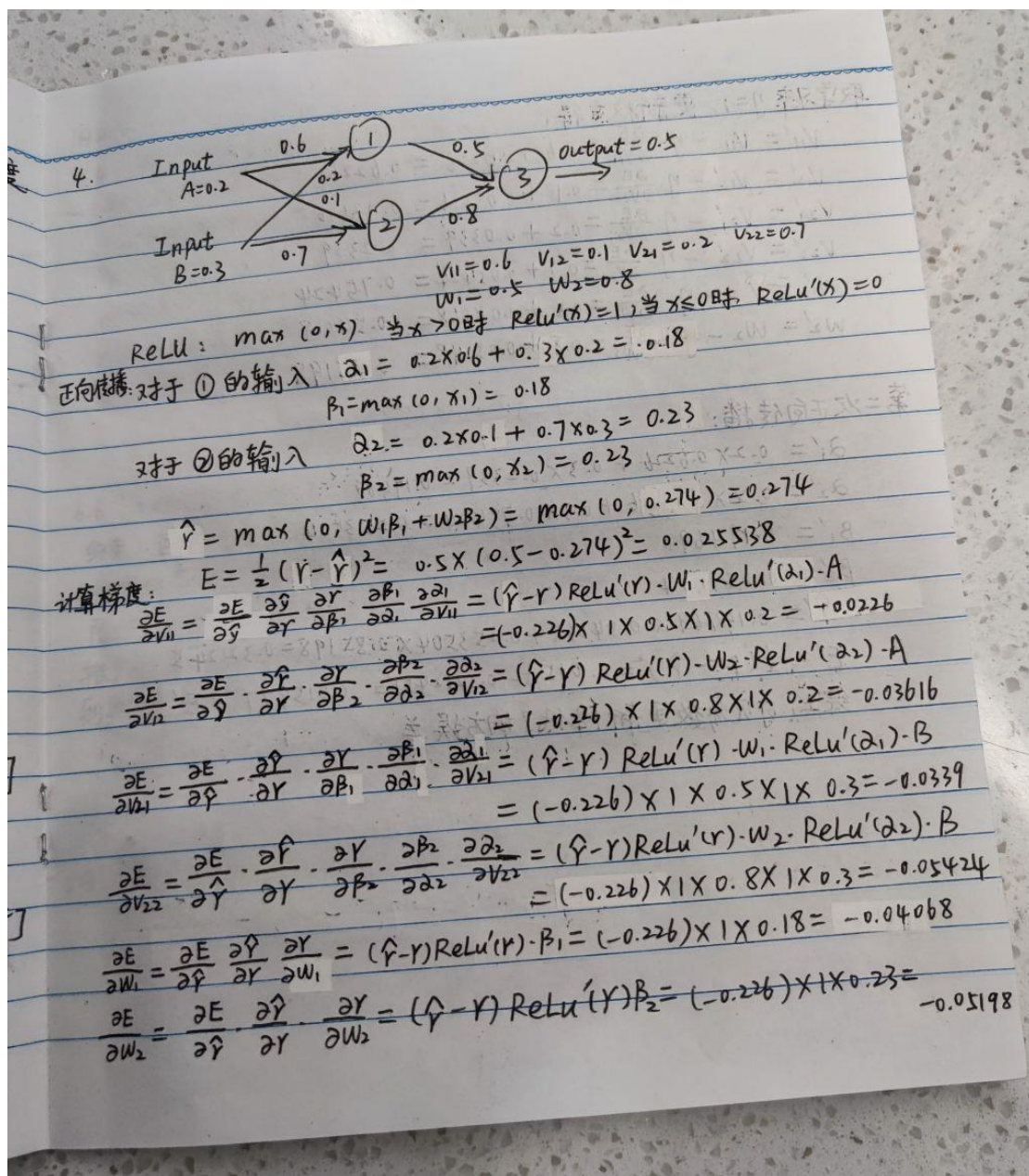
3. 计算 $\frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$ 和 $\log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}$ 关于向量 $x = [x_1, \dots, x_c]$ 的梯度 4.

$$\frac{\partial \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}}{\partial x_k} = \begin{cases} -\frac{\exp(x_i + x_k)}{\sum_{j=1}^c \exp(x_j)^2}, & k \neq i \\ \frac{\exp(x_k) \sum_{j=1, j \neq k}^c \exp(x_j)}{\sum_{j=1}^c \exp(x_j)^2}, & k = i \end{cases}$$

$$\frac{\partial \left[\log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)} \right]}{\partial x_k} = \begin{cases} -\frac{\exp(x_k)}{\sum_{j=1}^c \exp(x_j)}, & k \neq i \\ 1 - \frac{\exp(x_k)}{\sum_{j=1}^c \exp(x_j)}, & k = i \end{cases}$$

即有 $\frac{\partial \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)}}{\partial x} = \frac{e^{x_i}}{(\sum_{j=1}^c e^{x_j})^2} [-e^{x_1}, \dots, -e^{x_{i-1}}, \sum_{\substack{m=1 \\ m \neq i}}^c e^{x_m}, \dots, -e^{x_c}]$

$$\frac{\partial \left[\log \frac{\exp(x_i)}{\sum_{j=1}^c \exp(x_j)} \right]}{\partial x} = \frac{1}{\sum_{j=1}^c e^{x_j}} [-e^{x_1}, \dots, -e^{x_{i-1}}, \sum_{\substack{m=1 \\ m \neq i}}^c e^{x_m}, \dots, -e^{x_c}]$$



取学习率 $\eta=1$, 更新权重得:

$$V_{11}' = V_{11} - \eta \cdot \frac{\partial E}{\partial V_{11}} = 0.6 + 0.0226 = 0.6226$$

$$V_{12}' = V_{12} - \eta \cdot \frac{\partial E}{\partial V_{12}} = 0.1 + 0.03616 = 0.13616$$

$$V_{21}' = V_{21} - \eta \cdot \frac{\partial E}{\partial V_{21}} = 0.2 + 0.0339 = 0.2339$$

$$V_{22}' = V_{22} - \eta \cdot \frac{\partial E}{\partial V_{22}} = 0.7 + 0.05424 = 0.75424$$

$$W_{11}' = W_{11} - \eta \cdot \frac{\partial E}{\partial W_{11}} = 0.5 + 0.04068 = 0.54068$$

$$W_{21}' = W_{21} - \eta \cdot \frac{\partial E}{\partial W_{21}} = 0.8 + 0.05198 = 0.85198$$

第二次正向传播:

$$\alpha_1' = 0.2 \times 0.6226 + 0.3 \times 0.2339 = 0.19469$$

$$\alpha_2' = 0.2 \times 0.13616 + 0.3 \times 0.75424 = 0.253504$$

$$\beta_1' = 0.19469$$

$$\beta_2' = 0.253504$$

$$\hat{y}' = 0.19469 \times 0.54068 + 0.253504 \times 0.85198 = 0.321245$$

$$E' = \frac{1}{2} (y - \hat{y}')^2 = \frac{1}{2} (0.5 - 0.321245)^2 = 0.015977 < E$$

综上, 可见参数更新降低了均方误差

HW6.

6.4 线性判别分析能够解决 n 分类问题，而线性核 SVM 只能解决二分类问题。当线性判别分析的投影向量和线性核 SVM 的超平面向量垂直的时候，SVM 的最大间隔就是线性判别所要求的异类投影点间距，同时在这种情况下，线性判别分析的同类样例的投影点也会被这个超平面所划分一个，使其间隔较小。所以 (1) 线性判别分析求解出来的投影向量和线性核 SVM 求解出来的超平面向量垂直。(2) 数据集只有两类。(3) 数据集线性可分时，SVM 和 LDA 等价。

6.6 (1) SVM 的基本形态是一个硬间隔分类器，它要求所有样本都满足硬间隔约束，因此噪声很容易影响 SVM 的学习。(2) 存在噪声时，SVM 容易受噪声信息的影响，将训练得到的超平面向两个类间靠拢，导致训练的泛化能力降低。尤其是当噪声成为支持向量时，会直接影响整个超平面。(3) 当 SVM 推广到使用核函数时，会得到一个更复杂的模型，此时噪声也会一并被映射到更高维的特征，可能会对训练造成意想不到的结果。

综上，SVM 对噪声敏感。

6.9 对率回归的 L_2 正则化目标函数。

$$L(\beta) = \sum_{i=1}^m (-y_i \beta^T x_i + \ln(1 + e^{\beta^T x_i}))$$

$$F = L(\beta) + \frac{1}{2} \|\beta\|^2$$

由表示定理可知， $\beta = \sum_{i=1}^m \alpha_i \phi(x_i)$ ，带入上式可知。

$$\begin{aligned}
 F &= \sum_{i=1}^m [-y_i \beta^T \phi(x_i) + \ln(1 + e^{\beta^T \phi(x_i)})] + \frac{\lambda}{2} \|\beta\|^2 \\
 &= \sum_{i=1}^m [-y_i \sum_{j=1}^m \alpha_j \phi(x_i) \phi(x_j) + \ln(1 + e^{\sum_{j=1}^m \alpha_j \phi(x_i) \phi(x_j)})] + \\
 &\quad \frac{\lambda}{2} \|\sum_{j=1}^m \alpha_j \phi(x_j)\|^2 \\
 &= \sum_{i=1}^m [-y_i \sum_{j=1}^m \alpha_j k(x_i, x_j) + \ln(1 + e^{\sum_{j=1}^m \alpha_j k(x_i, x_j)})] + \frac{\lambda}{2} \|\sum_{j=1}^m \alpha_j k(x_j)\|^2
 \end{aligned}$$

目标函数为 $\min F$, 得到 L_2 正则化下的核岭回归:

$$\begin{aligned}
 \psi. \text{ 把 } \max_{\alpha, \hat{\alpha}} g(\alpha, \hat{\alpha}) &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) k(x_i, x_j) \\
 &\quad + \sum_{i=1}^m (y_i(\hat{\alpha}_i - \alpha_i) - \epsilon(\hat{\alpha}_i + \alpha_i)) \\
 \text{s.t. } C \geq \alpha, \hat{\alpha} \geq 0 \text{ and } \sum_{i=1}^m (\alpha_i - \hat{\alpha}_i) &= 0. \\
 \text{转化为类似标准型 } \max_{\alpha} g(\alpha) &= d^T v - \frac{1}{2} \alpha^T K \alpha \\
 \text{s.t. } C \geq \alpha \geq 0 \text{ and } \alpha^T u &= 0 \text{ 的形式.}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{1} \text{ 令 } \alpha &= [\alpha_1, \alpha_2, \dots, \alpha_m]^T, y = [y_1, y_2, \dots, y_m]^T. \\
 k_{ij} &= k(x_i, x_j) \quad \epsilon^* = [\epsilon, \epsilon, \dots, \epsilon]^T. \\
 \alpha^* &= [\alpha, \hat{\alpha}]^T, v = [-y - \epsilon^*, y - \epsilon^*]^T, K^* = \begin{bmatrix} K & -K \\ -K & K \end{bmatrix}. \\
 u &= [\underbrace{1, \dots, 1}_{m^T}, \underbrace{-1, \dots, -1}_{m^T}]^T
 \end{aligned}$$

$$\begin{aligned} \text{则 } \sum_{i=1}^m (y_i(\alpha_i - \hat{\alpha}_i) - \varepsilon(\alpha_i + \hat{\alpha}_i)) &= \sum_{i=1}^m (\alpha_i(y_i - \varepsilon)) + \hat{\alpha}_i(y_i - \varepsilon) \\ &= \alpha^T(-y - \varepsilon^*) + \hat{\alpha}^T(y - \varepsilon^*) \\ &= \alpha^{*T}V \end{aligned}$$

$$-\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) K(x_i, x_j)$$

$$\begin{aligned} &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i K(x_i, x_j) \alpha_j - \alpha_i K(x_i, x_j) \hat{\alpha}_j - \hat{\alpha}_i K(x_i, x_j) \alpha_j + \hat{\alpha}_i K(x_i, x_j) \hat{\alpha}_j) \\ &= -\frac{1}{2} (\alpha^T K \alpha - \alpha^T K \hat{\alpha} - \hat{\alpha}^T K \alpha + \hat{\alpha}^T K \hat{\alpha}) \\ &= -\frac{1}{2} \alpha^{*T} K \alpha^* \end{aligned}$$

$$\text{即有原 } \max_{\alpha, \hat{\alpha}} g(\alpha, \hat{\alpha}) = -\frac{1}{2} \alpha^{*T} K \alpha^* + \alpha^{*T} V$$

再看约束条件: 由 $C \geq \alpha, \hat{\alpha} \geq 0$ 且 $\alpha^* = [\alpha, \hat{\alpha}]^T$ 可知 $\alpha^* = \alpha_i$ 或 $\hat{\alpha}_i, 0 \leq \alpha^* \leq C$.

$$\text{由 } \sum_{i=1}^m (\alpha_i - \hat{\alpha}_i) = 0 \text{ 可知 } \alpha^{*T} u = 0.$$

综上, SVM 的对偶问题可转化为题意所示的标准形.

② 在软间隔 SVM 中, $v=1, u=y, K[i, j] = y_i y_j K(x_i, x_j)$
 $K(x_i, x_j) = \phi(x_i)^T \phi(x_j) = (x_i^T x_j)^2$, 求 $\phi(x_i)$ 表达式

设 x 是 m 维向量, 则

$$K(x_i, x_j) = (x_i^T x_j)^2 = \left(\sum_{u=1}^m x_{iu} x_{ju} \right) \left(\sum_{v=1}^m x_{iv} x_{jv} \right)$$

$$= \sum_{1 \leq u, v \leq m} x_{iu} x_{iv} x_{ju} x_{jv}$$

$$= [x_{i1} x_{i2} \dots x_{im}, \dots, x_{21} x_{22} \dots x_{2m}, \dots, x_{m1} x_{m2} \dots x_{mm}]^T$$

则 $\phi(x_i)$ 为一个 m^2 维向量, 每个分量为 $x_{iu} x_{iv}, 1 \leq u, v \leq m$

$$\phi(x_i) = [x_{i1} x_{i1}, x_{i1} x_{i2}, \dots, x_{i1} x_{im}, \dots, x_{i2} x_{i1}, \dots, x_{i2} x_{im}, \dots, x_{im} x_{im}]^T$$

各分量互不相同.