

① SVM Loss.

$$L = \sum_{i=1}^n L_i. \quad L_i = \sum_{j \neq y_i} \max(0, (Wx_i)_j - (Wx_i)_{y_i} + \Delta). \triangleq \sum_{j \neq y_i} \max(0, \text{gap}_i(j, y_i)).$$

$$\forall i \in \{1, \dots, n\}, \quad \frac{\partial L_i}{\partial w_{ke}} = \sum_{j \neq y_i} a_i^{(ke)}(j, y_i).$$

$$\text{其中 } a_i^{(ke)}(j, y_i) = \begin{cases} 1\{j=k\} x_{il} - 1\{y_i=k\} x_{il}, & \text{gap}_i(j, y_i) > 0 \\ 0, & \text{gap}_i(j, y_i) < 0 \end{cases}$$



$$\forall i \in \{1, \dots, n\}, \quad \begin{cases} k=y_i \text{ 时}, & \frac{\partial L_i}{\partial w_{ke}} = \sum_{j \neq y_i} -x_{il} 1\{\text{gap}_i(j, y_i) > 0\} = -x_{il} \sum_{j \neq y_i} 1\{\text{gap}_i(j, y_i) > 0\} \\ k \neq y_i \text{ 时}, & \frac{\partial L_i}{\partial w_{ke}} = \sum_{j \neq y_i} x_{il} 1\{\text{gap}_i(j, y_i) > 0\} 1\{j=k\} \\ & = x_{il} 1\{\text{gap}_i(j, y_i) > 0\}. \end{cases} \quad (*)$$

$$\text{记 } 1\{\text{gap}_i(j, y_i) > 0\} = m(i, j) \quad \text{则 } (*) \text{ 可表为 } \begin{cases} k=y_i, & \frac{\partial L_i}{\partial w_{ke}} = -x_{il} \sum_{j \neq y_i} m(i, j) \\ k \neq y_i, & \frac{\partial L_i}{\partial w_{ke}} = x_{il} m(i, k). \end{cases}$$

向量化, 抛弃 l

$$\begin{cases} k=y_i & \frac{\partial L_i}{\partial w_k} = -x_i \sum_{j \neq y_i} m(i, j) \\ k \neq y_i & \frac{\partial L_i}{\partial w_k} = x_i m(i, k) \end{cases}$$

$$\text{向量化, 构造矩阵 } C(i, k) = \begin{cases} -\sum_{j \neq y_i} m(i, j) & k=y_i \\ m(i, k) & k \neq y_i \end{cases}$$

$$\frac{\partial L}{\partial w_k} = x_i C(i, k) = C(i, k) x_i \quad \begin{matrix} \nearrow \text{数} \\ \searrow \text{向量} \end{matrix}$$

$$\left(\frac{\partial L_i}{\partial w_k} = c(i,k) x_i \right)$$

↓ 向量化, 抛弃 i

$$\frac{\partial L}{\partial w_k} = (c(1,k), \dots, c(n,k)) X$$

↓ 抛弃 k

$$\frac{\partial L}{\partial w} = C^T X.$$

② Softmax ~~Cross~~ (Cross-entropy loss)

$$L = \frac{1}{N} \sum_i L_i + \lambda \sum_k \sum_l w_{kl}^2, \quad L_i = -\log \left(\frac{e^{s_{yi}}}{\sum_j e^{s_j}} \right)$$

$$= -s_{yi} + \log \left(\sum_j e^{s_j} \right). \quad (*)$$

同理, 先分析 $\frac{\partial L_i}{\partial w_{kl}}$. 据 (*) 式可得 (*) 上的等式不要求, 改用 (*) 式)

$$\frac{\partial L_i}{\partial w_{kl}} = -x_{il} \mathbb{I}\{k=y_i\} + \frac{1}{\sum_j e^{s_j}} \cdot \sum_j e^{s_j} \mathbb{I}\{k=j\} x_{il}$$

($s = Wx_i$, 一般 $\mathbb{I}\{k=y_i\}$, $\mathbb{I}\{k=j\}$ 的出现是因为取了 w_{kl} 的 y 分量或 j 分量)

↓ 向量化, 抛弃 l

$$\frac{\partial L_i}{\partial w_k} = -x_i \mathbb{I}\{k=y_i\} + \frac{1}{\sum_j e^{s_j}} e^{s_k} \cdot x_i$$

$$\left(\frac{\partial L_i}{\partial w_k} = -x_i I\{k=y_i\} + \frac{1}{\sum_j e^{s_j}} e^{s_k} \cdot x_i \right)$$

↓

$$\frac{\partial L_i}{\partial w_k} = \begin{cases} k=y_i \text{ 时, } = [-1 + \frac{e^{s_k}}{\sum e^{s_j}}] x_i \\ k \neq y_i \text{ 时, } = [\frac{e^{s_k}}{\sum e^{s_j}}] \cdot x_i \end{cases}$$

↓ 向量化, 构造 $A(i, k) = \begin{cases} \frac{e^{s_k}}{\sum e^{s_j}} - 1, & k=y_i \\ \frac{e^{s_k}}{\sum e^{s_j}}, & k \neq y_i \end{cases}$

$$\frac{\partial L_i}{\partial w_k} = A(i, k) x_i$$

(A 的计算中涉及 $\exp()$, 为防止大数溢出, 可用上下同乘^正常数的技巧

$$\frac{c e^{s_k}}{c \sum e^{s_j}} = \frac{e^{s_k - c'}}{\sum e^{s_j - c'}}, \quad \text{令 } c' = \max_j \{s_j\}, \text{ 则指数部分 } s_j - c'$$

被中心化))

↓ 向量化, 扔掉 i

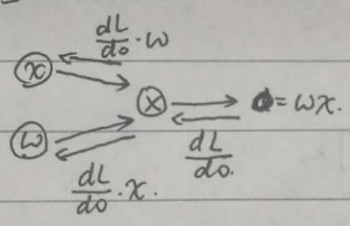
$$\frac{\partial L}{\partial w_k} = (A(1, k), \dots, A(n, k)) X.$$

↓ 扔掉 k.

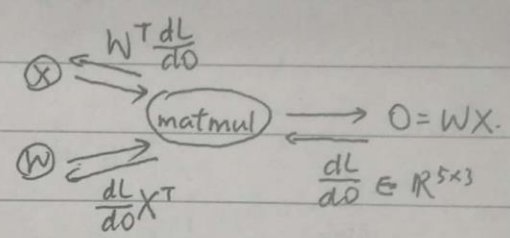
$$\frac{\partial L}{\partial w} = A^T X.$$

③ Two Layers NN

重要1:



向量化



($X \in \mathbb{R}^{10 \times 3}$, $W \in \mathbb{R}^{5 \times 10}$, $O \in \mathbb{R}^{5 \times 3}$)

(按照维度看是否转置)

(维度相同时套用此规律)

$$\left(\frac{\partial L}{\partial X} = \frac{\partial f}{\partial O} \cdot \frac{\partial O}{\partial X}, \quad \frac{\partial L}{\partial X_{kl}} = \sum_{ij} \frac{\partial L}{\partial O_{ij}} \cdot \frac{\partial O_{ij}}{\partial X_{kl}} \right. \quad (\text{因 } O = WX, \text{ 故 } j \neq l \text{ 时 } \frac{\partial O_{ij}}{\partial X_{kl}} = 0)$$

$$= \sum_i \frac{\partial L}{\partial O_{il}} \cdot \frac{\partial O_{il}}{\partial X_{kl}}$$

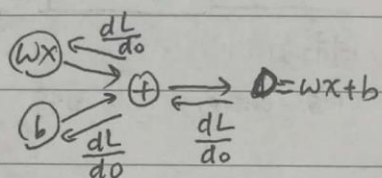
$$= \sum_i \frac{\partial L}{\partial O_{il}} \cdot w_{ik}$$

← 此处 l 是自由的! 该式意为
 $\frac{dL}{dO}$ 的第 l 列与 W 的第 k 列内积
 作为 $\frac{dL}{dX}$ 的第 k 行第 l 列元素

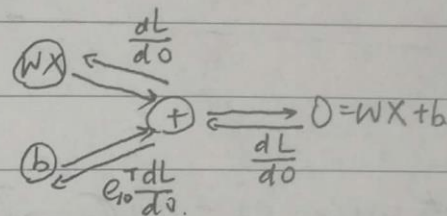
$$\frac{\partial L}{\partial X} = W^T \frac{dL}{dO}$$

关于 $\frac{\partial L}{\partial W}$ 可类似推导)

重要2:



向量化

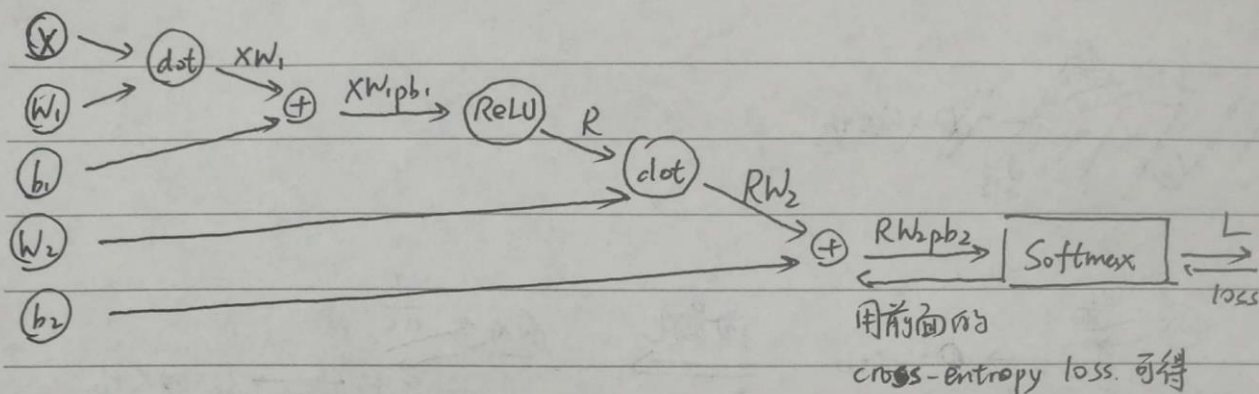


($WX \in \mathbb{R}^{10 \times 3}$, $b \in \mathbb{R}^{1 \times 3}$, $O \in \mathbb{R}^{10 \times 3}$)

(对于 broadcast 变量, 因最终的损失 L 是在行坐标上求和, 故梯度也应在行坐标上累加)

由上述两个重要性质, 结合下述 computation graph, 则可对 Two Layer NN 进行 BP.

-5-



用前面的 cross-entropy loss 可得

(若有正则化, 则应有 $W_1 \rightarrow \text{Softmax}$ 与 $W_2 \rightarrow \text{Softmax}$, 勿漏!)

-6-