

在 Softmax 回归中，我们利用交叉熵损失函数来计算训练集上的经验风险，即：

$$\mathcal{R}(W) = -\frac{1}{N} \sum_{n=1}^N (\mathbf{y}^{(n)})^T \log \hat{\mathbf{y}}^{(n)}$$

在训练过程中，需要计算风险函数对  $W$  的梯度，计算公式如下：

$$\frac{\partial \mathcal{R}(W)}{\partial W} = -\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} \left( \mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)} \right)^T$$

过程见下页↓

证: 设有  $K$  个类别, 则  $\hat{y}^{(n)}$  和  $y^{(n)}$  的形状为  $K \times 1$ ; 设  $x^{(n)}$  有  $D$  个特征, 则  $x^{(n)}$  的形状为  $D \times 1$ , 则  $W$  为形状为  $D \times K$  的矩阵.

令  $z^{(n)} = W^T \cdot x^{(n)} \Rightarrow z^{(n)}$  为  $K \times 1$  的列向量.

$$\text{则 } \hat{y}_i^{(n)} = \frac{\exp(z_i^{(n)})}{\sum_{j=1}^K \exp(z_j^{(n)})}$$

(1) 首先求 softmax 中  $\hat{y}_i^{(n)}$  对  $z_m^{(n)}$  的偏导.

① 当  $m=i$  时, 有

$$\begin{aligned} \frac{\partial \hat{y}_i^{(n)}}{\partial z_m^{(n)}} &= \frac{\partial \hat{y}_i^{(n)}}{\partial z_i^{(n)}} = \frac{\exp(z_i^{(n)}) \cdot \sum_{j=1}^K \exp(z_j^{(n)}) - \exp(z_i^{(n)}) \cdot \exp(z_i^{(n)})}{(\sum_{j=1}^K \exp(z_j^{(n)}))^2} \\ &= \hat{y}_i^{(n)} - (\hat{y}_i^{(n)})^2 \\ &= \hat{y}_i^{(n)} (1 - \hat{y}_i^{(n)}) \end{aligned}$$

② 当  $m \neq i$  时, 有

$$\frac{\partial \hat{y}_i^{(n)}}{\partial z_m^{(n)}} = \frac{-\exp(z_m^{(n)}) \cdot \exp(z_i^{(n)})}{(\sum_{j=1}^K \exp(z_j^{(n)}))^2} = -\hat{y}_m^{(n)} \hat{y}_i^{(n)}$$

(2) 接着求交叉熵损失对  $z_m^{(n)}$  的偏导.

$$\frac{\partial (y^{(n)T} \log \hat{y}^{(n)})}{\partial z_m^{(n)}} = \frac{\partial -(y^{(n)})^T \cdot \log \hat{y}^{(n)}}{\partial z_m^{(n)}} = \frac{\partial -\sum_{i=1}^K y_i^{(n)} \cdot \log \hat{y}_i^{(n)}}{\partial z_m^{(n)}}$$

$$= -\sum_{i=1}^K y_i^{(n)} \cdot \frac{1}{\hat{y}_i^{(n)}} \cdot \frac{\partial \hat{y}_i^{(n)}}{\partial z_m^{(n)}} \Rightarrow$$

$$\Rightarrow = -\left( y_m^{(n)} \cdot \frac{1}{\hat{y}_m^{(n)}} \cdot \frac{\partial \hat{y}_m^{(n)}}{\partial z_m^{(n)}} + \sum_{i \neq m} y_i^{(n)} \cdot \frac{1}{\hat{y}_i^{(n)}} \cdot \frac{\partial \hat{y}_i^{(n)}}{\partial z_m^{(n)}} \right) = -\left( y_m^{(n)} \cdot \frac{1 - \hat{y}_m^{(n)}}{\hat{y}_m^{(n)}} + \sum_{i \neq m} y_i^{(n)} \cdot \frac{-\hat{y}_m^{(n)} \hat{y}_i^{(n)}}{\hat{y}_i^{(n)}} \right)$$

$$= -\left[ y_m^{(n)} \cdot (1 - \hat{y}_m^{(n)}) + \sum_{i \neq m} y_i^{(n)} \cdot (-\hat{y}_m^{(n)}) \right]$$

1705546

$$= -\left[ y_m^{(n)} - y_m^{(n)} \cdot \hat{y}_m^{(n)} - \hat{y}_m^{(n)} \cdot \sum_{i \neq m} y_i^{(n)} \right] = \hat{y}_m^{(n)} - y_m^{(n)}$$



(3) 最后求交叉熵损失对  $w_{om}$  的偏导, 其中  $o=1,2,\dots,D, m=1,2,\dots,K$

$$\begin{aligned}\frac{\partial (-y^{(n)})^T \cdot \log \hat{y}^{(n)}}{\partial w_{om}} &= \frac{\partial (-y^{(n)})^T \cdot \log \hat{y}^{(n)}}{\partial z_m} \cdot \frac{\partial z_m}{\partial w_{om}} \\ &= (\hat{y}_m^{(n)} - y_m^{(n)}) \cdot \frac{\partial (\sum_{p=1}^D w_{mp}^{(n)} x_p^{(n)})}{\partial w_{om}} \\ &= (\hat{y}_m^{(n)} - y_m^{(n)}) \cdot x_o^{(n)}\end{aligned}$$

$$\Rightarrow \frac{\partial R(W)}{\partial W} = -\frac{1}{N} \sum_{n=1}^N \sum_{m=1}^K$$

$\Rightarrow$  对  $W$  的梯度如下所示:

$$\begin{bmatrix} (\hat{y}_1^{(n)} - y_1^{(n)}) \cdot x_1^{(n)} & \dots & (\hat{y}_K^{(n)} - y_K^{(n)}) \cdot x_1^{(n)} \\ \vdots & & \vdots \\ (\hat{y}_1^{(n)} - y_1^{(n)}) \cdot x_D^{(n)} & \dots & (\hat{y}_K^{(n)} - y_K^{(n)}) \cdot x_D^{(n)} \end{bmatrix}_{D \times K}$$

可写为  $x^{(n)} \cdot (\hat{y}^{(n)} - y^{(n)})^T$

对其中一个样本的第  $n$  个样本的交叉熵损失函数对  $W$  的梯度为  $x^{(n)} \cdot (\hat{y}^{(n)} - y^{(n)})^T$

$$\Rightarrow R \quad \frac{\partial R(W)}{\partial W} = -\frac{1}{N} \sum_{n=1}^N x^{(n)} \cdot (\hat{y}^{(n)} - y^{(n)})^T$$

$$\frac{\partial R(W)}{\partial W} = -\frac{1}{N} \sum_{n=1}^N x^{(n)} (\hat{y}^{(n)} - y^{(n)})^T$$

$$\text{即 } \frac{\partial R(W)}{\partial W} = -\frac{1}{N} \sum_{n=1}^N x^{(n)} (y^{(n)} - \hat{y}^{(n)})^T$$

