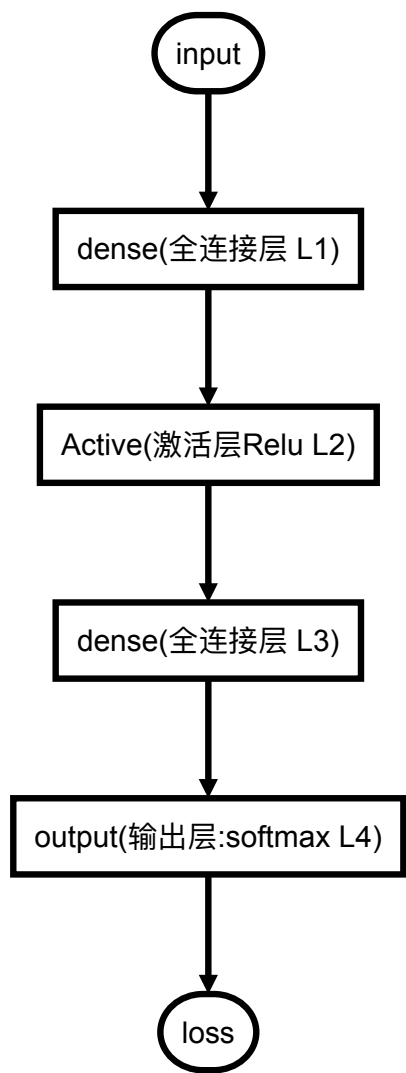


反向传播推倒

北京大学
地球与空间科学学院 陈宁 1601210216

本推倒考虑单个全连接隐藏层且激活函数为Relu，输出层为softmax，损失选择cross-entropy，为了灵活性，将全连接隐藏层拆分为两个层次——全连接dense层和激活层active层,即全连接层是一个线性变换不包含非线性因素。网络形式如下：



其中input为 $X^{(0)}$,经过 L_1 数据变为 $X^{(1)}$ ， 经过 L_2 数据变为 $X^{(2)}$ 经过 L_3 数据变为 $X^{(3)}$ 经过 L_4 数据变为 $X^{(4)}$ ， Forward过程为：

$$score = X^{(3)} = X^{(2)} * W_2 = Relu(X^{(1)}) * W_2 = Relu(X^{(0)} * W_1) * W_2$$

$$output = X^{(4)} = softmax(score)$$

$$loss = cross_entropy(output, Y)$$

反向传播过程：

(1) L_4 层 Loss的反向传播

$$\begin{aligned} L &= -\frac{1}{N} \sum_{i=1}^N \log\left(\frac{e^{x_{ik}}}{\sum_j e^{x_{ij}}}\right) \\ &= -\frac{1}{N} \sum_{i=1}^N [x_{ik} - \log \sum_j e^{x_{ij}}] \end{aligned}$$

其中N为batchsize,k为对应的true label, 因此梯度为

$$\begin{aligned} \frac{\partial L}{\partial X_{im}^{(3)}} &= \frac{1}{N} \sum_{i=1}^N \frac{e^{x_{im}}}{\sum_j e^{x_{ij}}} \quad (m \neq k) \\ \frac{\partial L}{\partial X_{ik}^{(3)}} &= \frac{1}{N} \sum_{i=1}^N \left[\frac{e^{x_{ik}}}{\sum_j e^{x_{ij}}} - 1 \right] \quad m = k \end{aligned}$$

(2) L_3 层反向传播过程:分别对W和X进行反向传播

$$\begin{aligned} dL &= tr\left[\left(\frac{\partial L}{\partial X^{(3)}}\right)^T \cdot d(X^{(2)} W_2)\right] \\ &= tr\left[\left(\frac{\partial L}{\partial X^{(3)}}\right)^T \cdot X^{(2)} dW_2\right] + tr\left[W_2 \left(\frac{\partial L}{\partial X^{(3)}}\right)^T dX^{(2)}\right] \end{aligned}$$

根据微分定义、标量对矩阵求导的定义对比可以得到：

$$\begin{aligned} dW_2 &= \frac{\partial L}{\partial W_2} = X^{(2)T} \cdot \frac{\partial L}{\partial X^{(3)}} \\ dX^{(2)} &= \frac{\partial L}{\partial X^{(2)}} = \frac{\partial L}{\partial X^{(3)}} \cdot W_2^T \end{aligned}$$

(3)对激活层反向传播过程 $Relu = \max(0, x)$

$$\frac{\partial L}{\partial X_{ij}^{(1)}} = \begin{cases} dX_{ij}^{(2)} \cdot 1 & \text{对于 } X_{ij}^{(1)} > 0 \text{ 的部分特征} \\ 0 & \text{otherwise} \end{cases}$$

(4)对 L_1 全连接层的反向传播只需要对 W_1 梯度

$$\frac{\partial L}{\partial W_1} = X^{(0)T} \bullet dX^{(1)}$$