

神经网络的梯度反向传播的数学推导

全连接层

简单全连接网络的反向传播过程，
设第 l 层，

$$z^l = W^l \cdot x^{l-1} + b^l$$

$$x^{l-1} \in R^m$$

$$x^l = f(z^l)$$

$$x^l \in R^n$$

其中 f 是非线性激励函数，

$$W^l \in R^{n \times m}$$

$J = J(W, b)$ 是损失函数， x^{l-1} 是第 $l-1$ 层的输出。

$$\text{记 } \delta^l = \frac{\partial J}{\partial z^l}, \quad \therefore z_i^l = \sum_j w_{ij}^l \cdot x_j^{l-1} + b_i^l \quad (i=0, 1, 2, \dots, n-1)$$

$$\therefore \frac{\partial J}{\partial w_{ij}^l} = \frac{\partial J}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial w_{ij}^l} = \delta_i^l \cdot x_j^{l-1}$$

$$\Rightarrow \boxed{\frac{\partial J}{\partial W^l} = \delta^l \cdot (x^{l-1})^T} \quad ①$$

$$\frac{\partial J}{\partial b_i^l} = \frac{\partial J}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial b_i^l}$$

$$\Rightarrow \boxed{\frac{\partial J}{\partial b^l} = \delta^l} \quad ②$$

$$\therefore z^l = W^l \cdot f(z^{l-1}) + b^l$$

$$\therefore z_i^l = \sum_j w_{ij}^l \cdot f(z_j^{l-1}) + b_i^l \quad (i=0, 1, \dots, n-1)$$

$$\therefore \frac{\partial J}{\partial z_i^{l-1}} = \sum_k \frac{\partial J}{\partial z_k^l} \cdot \frac{\partial z_k^l}{\partial z_i^{l-1}}$$

$$= \sum_k \delta_k^l \cdot w_{ki}^l \cdot f'(z_i^{l-1})$$

$$\Rightarrow \delta_i^{l-1} = \left(\sum_k \delta_k^l \cdot w_{ki}^l \right) \cdot f'(z_i^{l-1}) \Rightarrow \boxed{\delta^{l-1} = (W^l)^T \cdot \delta^l \odot f'(z^{l-1})} \quad ③$$

$$\text{对于输出层 } L \text{ 而言: } \delta^L = \frac{\partial J}{\partial z^L} = \frac{\partial J}{\partial x^L} \odot \frac{\partial x^L}{\partial z^L}$$

$$\Rightarrow \boxed{\delta^L = \frac{\partial J}{\partial x^L} \odot f'(z^L)} \quad ④$$

③ 是矩阵 Hadamard 积，
也就是逐元素相乘。

最后总结为：

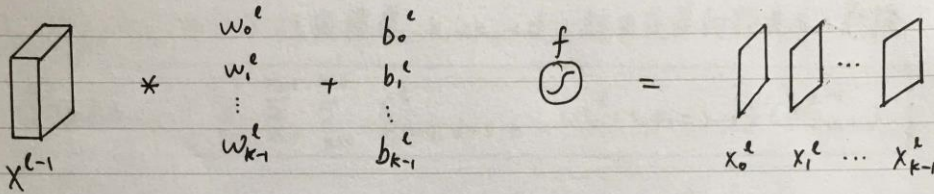
- 1) 由训练样本给网络输入 X 和输出 Y 。
- 2) 从第 2 层到第 L 层进行前向网络计算
- 3) 由公式④ 计算第 L 层的残差 δ^L
- 4) 从第 L 层 到第 2 层，由公式③ 逐层计算各层残差 δ^l ，然后由公式①、② 计算梯度

卷积层

卷积层的误差反向传播 [参考 3)4)5)6)]

设第 l 层为卷积层. 输入 $X^{l-1} \in \mathbb{R}^{s_1 \times s_2 \times d}$, 有 k 个卷积核 $w_0^l, w_1^l, \dots, w_{k-1}^l$, 且 $w_k^l \in \mathbb{R}^{k_1 \times k_2 \times d}$, 输出为 $X^l \in \mathbb{R}^{q_1 \times q_2 \times k}$ 其中 $q_1 = s_1 - k_1 + 1, q_2 = s_2 - k_2 + 1$ (以 stride=1 为例)

记: $W^l = [w_0^l, w_1^l, \dots, w_{k-1}^l]$, 每一个卷积核 w_k^l 对应一个实数偏置 $b_k^l \in \mathbb{R}$



x_k^l 是第 k 个卷积核对应的特征图 (feature map), 这 k 个特征图按通道叠起来就是输出 X^l 了.

$$z^l = X^{l-1} * W^l + b^l \quad \text{其中 } z^l \in \mathbb{R}^{q_1 \times q_2 \times k} \text{ 为中间值,}$$

$$X^l = f(z^l) \quad \text{f 是激活函数.}$$

设 J 是网络损失函数. 记 $\frac{\partial J}{\partial z^l} = \delta^l$ 为误差项

由 $z^l = f(z^{l-1}) * W^l + b^l$, 按卷积定义:

$$z_{i,j,k}^l = \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} \sum_{c=0}^{d-1} f(z_{i+m,j+n,c}^{l-1}) \cdot w_k^l(m,n,c) + b_k^l$$

$$(0 \leq i < q_1, \quad 0 \leq j < q_2)$$

$$\frac{\partial J}{\partial z_{u,v,c}^{l-1}} = \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \sum_{k=0}^{k-1} \frac{\partial J}{\partial z_{i,j,k}^l} \cdot \frac{\partial z_{i,j,k}^l}{\partial z_{u,v,c}^{l-1}} \quad \begin{matrix} 0 \leq u < s_1 \\ 0 \leq v < s_2 \\ 0 \leq c < d \end{matrix}$$

$$\text{即: } \delta_{u,v,c}^{l-1} = \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \sum_{k=0}^{k-1} \delta_{i,j,k}^l \cdot \frac{\partial z_{i,j,k}^l}{\partial z_{u,v,c}^{l-1}}$$

$$\text{由 } \frac{\partial z_{i,j,k}^l}{\partial z_{u,v,c}^{l-1}} = f'(z_{u,v,c}^{l-1}) \cdot w_k^l(u-i, v-j, c)$$

$$\therefore \delta_{u,v,c}^{l-1} = \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \sum_{k=0}^{K-1} \delta_{i,j,k}^l \cdot w_k^l(u-i, v-j, c) \cdot f'(z_{u,v,c}^{l-1})$$

对于 w_k^l 而言下标范围是 $k_1 \times k_2 \times d$, 超出范围的都是0. 更换下标和下标

$$\text{得: } \delta_{u,v,c}^{l-1} = \sum_{s=0}^{k_1-1} \sum_{t=0}^{k_2-1} \sum_{k=0}^{K-1} \delta_{u-s, v-t, k}^l \cdot w_k^l(s, t, c) \cdot f'(z_{u,v,c}^{l-1}) \quad ①$$

这就是卷积层的误差反向传播公式, 其中的求和项.

$\delta_{u-s, v-t, k}^l \cdot w_k^l(s, t, c)$, 不考虑通道的话, 其实就是把 w_k^l 进行上下翻转, 再进行左右翻转后作为卷积核与 δ^l 进行卷积.

有了 δ 后再计算梯度, $\frac{\partial J}{\partial w}$, $\frac{\partial J}{\partial b}$

$$\therefore z_{i,j,k}^l = \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} \sum_c X_{i+m, j+n, c}^{l-1} \cdot w_k^l(m, n, c) + b_k^l$$

由于第 k 个卷积核只影响到第 k 个特征图, 所以:

$$\begin{aligned} \frac{\partial J}{\partial w_k^l(m, n, c)} &= \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \frac{\partial J}{\partial z_{i,j,k}^l} \cdot \frac{\partial z_{i,j,k}^l}{\partial w_k^l(m, n, c)} \\ &= \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \delta_{i,j,k}^l \cdot X_{i+m, j+n, c}^{l-1} \end{aligned} \quad \text{可以看出这也是一个卷积计算.}$$

$$\text{也即: } \frac{\partial J}{\partial w_k^l(c)} = X_{(c)}^{l-1} * \delta_k^l$$

\downarrow
 $R^{k_1 \times k_2}$

\downarrow
 $R^{s_1 \times s_2}$

\downarrow
 $R^{q_1 \times q_2}$

②

$$\begin{aligned} \frac{\partial J}{\partial b_k^l} &= \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \frac{\partial J}{\partial z_{i,j,k}^l} \cdot \frac{\partial z_{i,j,k}^l}{\partial b_k^l} \\ &= \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \delta_{i,j,k}^l \end{aligned}$$

$$\Rightarrow \frac{\partial J}{\partial b_k^l} = \sum_{i=0}^{q_1-1} \sum_{j=0}^{q_2-1} \delta_{i,j,k}^l \quad ③$$

与全连接层类似: 由式子 (1) 从输出层开始反向递推计算各层的误差项, 然后由式子 (2) (3) 分别计算各层的梯度 (损失函数对卷积核和偏置的导数)

池化层[参考 8]

平均池化

以一维池化核为 2 的为例。

$$X_i^l = \frac{1}{2} \sum_{u=i}^{i+1} X_u^{l-1}$$
$$\delta_u^{l-1} = \frac{\partial J}{\partial X_u^{l-1}} = \sum_i \frac{\partial J}{\partial X_i^l} \cdot \frac{\partial X_i^l}{\partial X_u^{l-1}} \quad \begin{array}{l} \because u = i, i+1 \\ \therefore i \text{ 取值 } u, u-1 \end{array}$$
$$= \sum_{i=u-1}^u \frac{\partial J}{\partial X_i^l} \cdot \frac{\partial X_i^l}{\partial X_u^{l-1}}$$
$$\Rightarrow \delta_u^{l-1} = \sum_{i=u-1}^u \delta_i^l \cdot \frac{1}{2}$$
$$= \frac{1}{2} \sum_{i=u-1}^u \delta_i^l$$

最大池化

$$X_i^l = \max_{u=i, \dots, i+1} \{ X_u^{l-1} \}$$
$$\frac{\partial J}{\partial X_u^{l-1}} = \sum_i \frac{\partial J}{\partial X_i^l} \cdot \frac{\partial X_i^l}{\partial X_u^{l-1}} \quad i = u-1, u$$
$$\Rightarrow \delta_u^{l-1} = \sum_{i=u-1}^u \delta_i^l \cdot \frac{\partial X_i^l}{\partial X_u^{l-1}}$$
$$\text{而 } \frac{\partial X_i^l}{\partial X_u^{l-1}} = \begin{cases} 1 & X_i^l \text{ 取最大值是来自于 } X_u^{l-1}, \text{ 则为 } 1, \text{ 否则为 } 0 \\ 0 \end{cases}$$

[参考文献]

- 1) [Neural Networks and Deep Learning](<http://neuralnetworksanddeeplearning.com/>) by Michael Nielsen</br>
- 2) [Deep Learning](<http://www.deeplearningbook.org/>), book by Ian Goodfellow, Yoshua Bengio, and Aaron Courville</br>
- 3) <http://cs231n.github.io/convolutional-networks> </br>
- 4) http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture05.pdf
- 5) <http://ufldl.stanford.edu/tutorial/supervised/ConvolutionalNeuralNetwork/>
- 6) http://deeplearning.net/software/theano_versions/dev/tutorial/conv_arithmetic.html#transposed-convolution-arithmetic
- 7) http://ufldl.stanford.edu/wiki/index.php/Gradient_checking_and_advanced_optimization
- 8) https://software.intel.com/sites/products/documentation/doclib/daal/daal-user-and-reference-guides/daal_prog_guide/GUID-2C3AA967-AE6A-4162-84EB-93BE438E3A05.htm