

# Feedforward NN 与 Convolutional NN 中的 Back Propagation

## 1 引入

经过 forward propagation 之后，我们可以得到当前神经网络的 cost function  $J(\mathbf{W}, \mathbf{b})$ ，为了找到 cost function 的最小值，我们需要沿其梯度方向更新每层网络的权重矩阵  $\mathbf{W}$  和偏置列向量  $\mathbf{b}$  的值：

$$\mathbf{W}^{[l]} := \mathbf{W}^{[l]} - \frac{\partial J}{\partial \mathbf{W}^{[l]}} \quad (1)$$

$$\mathbf{b}^{[l]} := \mathbf{b}^{[l]} - \frac{\partial J}{\partial \mathbf{b}^{[l]}} \quad (2)$$

所以我们需要计算  $\frac{\partial J}{\partial \mathbf{W}^{[l]}}$  和  $\frac{\partial J}{\partial \mathbf{b}^{[l]}}$ 。

由于直接求数值积分的高计算消耗，我们利用 computation graph 的逻辑，利用链式法则逐层迭代地对  $\nabla_{\mathbf{W}} J$  和  $\nabla_{\mathbf{b}} J$  进行计算。

## 2 Feedforward NN 的 Back Propagation

由于

$$\frac{\partial J}{\partial \mathbf{W}^{[l]}} = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{W}^{[l]}} \quad (3)$$

所以我们想要先计算  $\frac{\partial J}{\partial \mathbf{Z}^{[l]}}$ ，再利用链式法则计算  $\frac{\partial J}{\partial \mathbf{W}^{[l]}}$ 。

又由于

$$\frac{\partial J}{\partial \mathbf{Z}^{[l]}} = \frac{\partial J}{\partial \mathbf{A}^{[l]}} \frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}} \quad (4)$$

所以我们最终发现应该先计算  $\frac{\partial J}{\partial \mathbf{A}^{[l]}}$ 。

## 2.1 计算 $\partial J / \partial \mathbf{A}^{[l]}$

对于  $J$  的全微分, 有:

$$dJ = \sum_{i=1}^p \sum_{j=1}^q \frac{\partial J}{\partial Z_{ij}^{[l]}} dZ_{ij}^{[l]} = \sum_{i,j} (\nabla_{Z^{[l]}} J)_{ij} dZ_{ij}^{[l]} = \langle \nabla_{Z^{[l]}} J, dZ^{[l]} \rangle_F = \text{tr} \left( (\nabla_{Z^{[l]}} J)^\top dZ^{[l]} \right) \quad (5)$$

代入  $dZ^{[l]} = \mathbf{W}^{[l]} d\mathbf{A}^{[l-1]}$ , 得:

$$dJ = \text{tr} \left( (\nabla_{Z^{[l]}} J)^\top \mathbf{W}^{[l]} d\mathbf{A}^{[l-1]} \right) \quad (6)$$

所以:

$$\nabla_{\mathbf{A}^{[l-1]}} J = \left( (\nabla_{Z^{[l]}} J)^\top \mathbf{W}^{[l]} \right)^\top \quad (7)$$

由于  $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$ , 所以:

$$\frac{\partial J}{\partial \mathbf{A}^{[l-1]}} = \mathbf{W}^{[l]\top} \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \quad (8)$$

从  $\mathbf{A}^{[l-1]}$  推导  $\mathbf{A}^{[l]}$ :

$$\frac{\partial J}{\partial \mathbf{A}^{[l]}} = \mathbf{W}^{[l+1]\top} \frac{\partial J}{\partial \mathbf{Z}^{[l+1]}} \quad (9)$$

## 2.2 计算 $\partial J / \partial \mathbf{Z}^{[l]}$

和 (4) 式相同地, 有:

$$\frac{\partial J}{\partial \mathbf{Z}^{[l]}} = \frac{\partial J}{\partial \mathbf{A}^{[l]}} \frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}}$$

我们刚才已经求得了  $\frac{\partial J}{\partial \mathbf{A}^{[l]}}$ 。这里我们分析一下  $\frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}}$ 。由于 activation function element-wise 地作用在  $\mathbf{Z}$ , 所以  $\frac{\partial \mathbf{A}}{\partial \mathbf{Z}}$  得到的 Jacobian 为对角矩阵。因此可以用 Hadamard 乘 (也就是 element-wise 乘法) 实现同样效果, 即:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{Z}^{[l]}} &= \frac{\partial J}{\partial \mathbf{A}^{[l]}} \mathbf{J}_g = \frac{\partial J}{\partial \mathbf{A}^{[l]}} \text{diag} \left( g^{[l]'}(\mathbf{Z}^{[l]}) \right) \\ &= \frac{\partial J}{\partial \mathbf{A}^{[l]}} \odot g^{[l]'}(\mathbf{Z}^{[l]}) \end{aligned} \quad (10)$$

$$\left( = \mathbf{W}^{[l+1]\top} \frac{\partial J}{\partial \mathbf{Z}^{[l+1]}} \odot g^{[l]'}(\mathbf{Z}^{[l]}) \right) \quad (11)$$

## 2.3 计算 $\partial J / \partial \mathbf{W}^{[l]}$

和 (3) 式相同地, 有:

$$\frac{\partial J}{\partial \mathbf{W}^{[l]}} = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{W}^{[l]}}$$

我们刚才已经求得了  $\frac{\partial J}{\partial \mathbf{Z}^{[l]}}$ 。这里我们可以接着用 vec-Kronecker 恒等式求  $\frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{W}^{[l]}}$ , 或者模仿  $\frac{\partial J}{\partial \mathbf{A}^{[l]}}$  的推导, 从全微分和迹运算进行推导 (相对简单)。

### 2.3.1 硬求 $\partial \mathbf{Z}^{[l]}/\partial \mathbf{W}^{[l]}$

对式  $d\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} d\mathbf{A}^{[l-1]}$  的两侧同时做  $\text{vec}$  运算并插入单位矩阵, 根据恒等式  $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B})$  得到:

$$\text{vec}(d\mathbf{Z}^{[l]}) = (\mathbf{A}^{[l-1]\top} \otimes I_{n_l}) \text{vec}(d\mathbf{W}^{[l]}) \quad (12)$$

其中“ $\otimes$ ”表示 Kronecker 积。所以  $\frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{W}^{[l]}}$  得到的 Jacobian (vec 形式) 为:

$$\frac{\partial \text{vec}(\mathbf{Z}^{[l]})}{\partial \text{vec}(\mathbf{W}^{[l]})} = \mathbf{A}^{[l-1]\top} \otimes I_{n_l} \quad (13)$$

同样改写  $\frac{\partial J}{\partial \mathbf{W}^{[l]}}$  为 vec 形式, 有:

$$\frac{\partial J}{\partial \text{vec}(\mathbf{W}^{[l]})} = \left( \frac{\partial \text{vec}(\mathbf{Z}^{[l]})}{\partial \text{vec}(\mathbf{W}^{[l]})} \right)^\top \frac{\partial J}{\partial \text{vec}(\mathbf{Z}^{[l]})} \quad (14)$$

将式 (12) 代入式 (13), 有:

$$\text{vec}\left(\frac{\partial J}{\partial \mathbf{W}^{[l]}}\right) = (\mathbf{A}^{[l-1]\top} \otimes I_{n_l})^\top \text{vec}\left(\frac{\partial J}{\partial \mathbf{Z}^{[l]}}\right) \quad (15)$$

由于  $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$ , 所以:

$$\text{vec}\left(\frac{\partial J}{\partial \mathbf{W}^{[l]}}\right) = (\mathbf{A}^{[l-1]} \otimes I_{n_l}) \text{vec}\left(\frac{\partial J}{\partial \mathbf{Z}^{[l]}}\right) \quad (16)$$

根据恒等式  $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B})$  得到:

$$\text{vec}\left(\frac{\partial J}{\partial \mathbf{W}^{[l]}}\right) = I_{n_l} \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \mathbf{A}^{[l-1]\top} \quad (17)$$

因为  $I_{n_l}$  左乘并不改变数值, 所以:

$$\text{vec}\left(\frac{\partial J}{\partial \mathbf{W}^{[l]}}\right) = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \mathbf{A}^{[l-1]\top} \Rightarrow \frac{\partial J}{\partial \mathbf{W}^{[l]}} = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \mathbf{A}^{[l-1]\top} \quad (18)$$

### 2.3.2 $J$ 的全微分

由  $J$  的全积分

$$dJ = \text{tr}\left((\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{Z}^{[l]}\right)$$

代入  $d\mathbf{Z}^{[l]} = d\mathbf{W}^{[l]} \mathbf{A}^{[l-1]}$ , 得:

$$dJ = \text{tr}\left((\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{W}^{[l]} \mathbf{A}^{[l-1]}\right)$$

由迹运算的循环不变性即  $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{CAB})$  得到:

$$\text{tr}\left((\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{W}^{[l]} \mathbf{A}^{[l-1]}\right) = \text{tr}\left(\mathbf{A}^{[l-1]} (\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{W}^{[l]}\right) = \text{tr}\left((\nabla_{\mathbf{Z}^{[l]}} J \mathbf{A}^{[l-1]\top})^\top d\mathbf{W}^{[l]}\right) \quad (19)$$

所以:

$$\frac{\partial J}{\partial \mathbf{W}^{[l]}} = \nabla_{\mathbf{Z}^{[l]}} J \mathbf{A}^{[l-1]\top} = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \mathbf{A}^{[l-1]\top} \quad (20)$$

## 2.4 计算 $\partial J / \partial \mathbf{b}^{[l]}$

和计算  $\frac{\partial J}{\partial \mathbf{W}^{[l]}}$  的过程类似地，我们也有两种方法计算  $\frac{\partial J}{\partial \mathbf{b}^{[l]}}$ 。这里我们只介绍第二种迹运算的方法。第一种利用 **vec-Kronecker** 恒等式的方法留给读者作为练习。

由 forward propagation 的公式：

$$\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} \mathbf{A}^{[l-1]} + \mathbf{b}^{[l]} \mathbf{1}_{1 \times m}, \mathbf{b} \in \mathbb{R}^{n_l \times 1}$$

可以得到：

$$d\mathbf{Z}^{[l]} = d(\mathbf{b}^{[l]} \mathbf{1}_{1 \times m}) = (d\mathbf{b}^{[l]}) \mathbf{1}_{1 \times m} + \mathbf{b}^{[l]} \underbrace{(d\mathbf{1}_{1 \times m})}_{=0} = d\mathbf{b}^{[l]} \mathbf{1}_{1 \times m} \quad (21)$$

由  $J$  的全积分

$$dJ = \text{tr} \left( (\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{Z}^{[l]} \right)$$

代入  $d\mathbf{Z}^{[l]} = d\mathbf{b}^{[l]} \mathbf{1}_{1 \times m}$  (其中  $m$  为样本数)，得：

$$\text{tr} \left( (\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{Z}^{[l]} \right) = \text{tr} \left( (\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{b}^{[l]} \mathbf{1}_{1 \times m} \right) \quad (22)$$

由迹运算的循环不变性，得：

$$\text{tr} \left( (\nabla_{\mathbf{Z}^{[l]}} J)^\top d\mathbf{b}^{[l]} \mathbf{1}_{1 \times m} \right) = \text{tr} \left( \left( \nabla_{\mathbf{Z}^{[l]}} J (\mathbf{1}_{1 \times m})^\top \right)^\top d\mathbf{b}^{[l]} \right) = \text{tr} \left( (\nabla_{\mathbf{Z}^{[l]}} J \mathbf{1}_{m \times 1})^\top d\mathbf{b}^{[l]} \right) \quad (23)$$

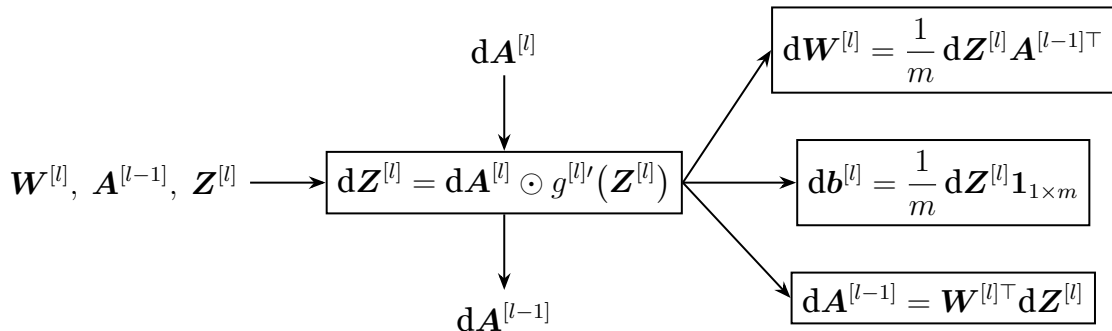
所以：

$$\frac{\partial J}{\partial \mathbf{b}^{[l]}} = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \mathbf{1}_{m \times 1} \quad (24)$$

这即等价于对  $\frac{\partial J}{\partial \mathbf{Z}^{[l]}}$  逐行求和。故在程序中写作 `np.sum(dZ, axis=1, keepdims=True)`。

## 2.5 总结：Feedforward NN 的 Back Propagation

整合之前的推导，我们可以得到：



### 3 Convolutional Layer 的 Back Propagation

注意：

1. 我们用 “\*” 表示 convolution，用 “ $\star$ ” 表示 cross-correlation。
2. 本节不会在推导中大量使用矩阵运算，因为卷积操作可以用行列索引方便表示。
3. 下标  $[a, b]$  表示第  $a$  行第  $b$  列的元素，索引从  $[0, 0]$  开始。 $f^{[l]}$  表示  $l$  层的 kernel/filter 的行/列数。 $l$  层的激活张量  $\mathbf{A}^{[l]}$  的 shape 为  $[n_H^{[l-1]}, n_W^{[l-1]}, n_c^{[l-1]}, m]$ 。其中  $m$  为样本数， $n_c^{[l-1]}$  为  $l-1$  层的 channel 数。
4. 阅读的时候可以认为每层神经网络的 kernel/filter 的数量为 1，且该 kernel/filter 的层数也为 1。

#### 3.1 从 $\partial J / \partial \mathbf{Z}^{[l]}$ 推导 $\partial J / \partial \mathbf{W}^{[l]}$

由链式法则：

$$\frac{\partial J}{\partial w_{[a,b]}^{[l]}} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} \frac{\partial z_{[x,y]}^{[l]}}{\partial w_{[a,b]}^{[l]}} \quad (25)$$

其中， $z_{[x,y]}^{[l]}$  可以表示为：

$$z_{[x,y]}^{[l]} = a^{[l-1]} \star w^{[l]} + b^{[l]} = \sum_{a=0}^{f^{[l]}} \sum_{b=0}^{f^{[l]}} a_{[x+a,y+b]}^{[l-1]} w_{[a,b]}^{[l]} + b^{[l]} \quad (26)$$

即对应范围的卷积再加上偏置项（在 implement forward propagation from scratch 的时候其实也是将对应的区域 slice 出来再进行卷积操作）。

将 (26) 式代入 (25) 式得到：

$$\frac{\partial J}{\partial w_{[a,b]}^{[l]}} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} \frac{\partial \left( \sum_{a'=0}^{f^{[l]}} \sum_{b'=0}^{f^{[l]}} a_{[x+a',y+b']}^{[l-1]} w_{[a',b']}^{[l]} + b^{[l]} \right)}{\partial w_{[a,b]}^{[l]}} \quad (27)$$

注意！当且仅当  $a' = a$  且  $b' = b$  的时候第二项偏导不为 0。所以：

$$\frac{\partial J}{\partial w_{[a,b]}^{[l]}} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} a_{[x+a,y+b]}^{[l-1]} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} g \left( z_{[x+a,y+b]}^{[l-1]} \right) \quad (28)$$

所以：

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{W}^{[l]}} &= \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \star g(\mathbf{Z}^{[l-1]}) \\ &\left( = \frac{\partial J}{\partial \mathbf{Z}^{[l]}} * g(\text{ROT180}(\mathbf{Z}^{[l-1]})) \right) \end{aligned} \quad (29)$$

其中 ROT180 表示上下翻转一次，接着左右翻转一次。

### 3.2 从 $\partial J/\partial \mathbf{Z}^{[l]}$ 推导 $\partial J/\partial \mathbf{b}^{[l]}$

与式 (27) 类似地：

$$\frac{\partial J}{\partial b^{[l]}} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} \frac{\partial \left( \sum_{a'=0}^{f^{[l]}} \sum_{b'=0}^{f^{[l]}} a_{[x+a',y+b']}^{[l-1]} w_{[a',b']}^{[l]} + b^{[l]} \right)}{\partial b^{[l]}} \quad (30)$$

显然第二项偏导等于 1。所以：

$$\frac{\partial J}{\partial b^{[l]}} = \sum_{x=0}^{n_H^{[l]}} \sum_{y=0}^{n_W^{[l]}} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} = \mathbf{1}_{n_H^{[l]} \times n_W^{[l]}} \star \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \quad (31)$$

### 3.3 推导 $\partial J/\partial \mathbf{Z}^{[l]}$

综合式 (29) 与式 (31)，可以得到：

$$\begin{aligned} \frac{\partial J}{\partial w^{[l]}} &= g(\mathbf{Z}^{[l-1]}) \star \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \\ \frac{\partial J}{\partial b^{[l]}} &= \mathbf{1}_{n_H^{[l]} \times n_W^{[l]}} \star \frac{\partial J}{\partial \mathbf{Z}^{[l]}} \end{aligned}$$

所以我们要求得  $\frac{\partial J}{\partial \mathbf{Z}^{[l]}}$ 。

显然地，我们有：

$$z_{[x',y']}^{[l+1]} = \sum_{a=0}^{f^{[l+1]}} \sum_{b=0}^{f^{[l+1]}} \left( g \left( z_{[x'+a,y'+b]}^{[l]} \right) \right) w_{[a,b]}^{[l+1]} + b^{[l+1]} \quad (32)$$

$$\frac{\partial J}{\partial z_{[x,y]}^{[l]}} = \sum_{x'=0}^{n_H^{[l+1]}} \sum_{y'=0}^{n_W^{[l+1]}} \frac{\partial J}{\partial z_{[x',y']}^{[l+1]}} \frac{\partial z_{[x',y']}^{[l+1]}}{\partial z_{[x,y]}^{[l]}} \quad (33)$$

将式 (32) 代入式 (33)，得到：

$$\frac{\partial J}{\partial z_{[x,y]}^{[l]}} = \sum_{x'=0}^{n_H^{[l+1]}} \sum_{y'=0}^{n_W^{[l+1]}} \frac{\partial J}{\partial z_{[x',y']}^{[l+1]}} \frac{\partial \sum_{a=0}^{f^{[l+1]}} \sum_{b=0}^{f^{[l+1]}} \left( g \left( z_{[x'+a,y'+b]}^{[l]} \right) \right) w_{[a,b]}^{[l+1]} + b^{[l+1]}}{\partial z_{[x,y]}^{[l]}} \quad (34)$$

同样地，当且仅当  $x' + a = x$  且  $y' + b = y$  时第二项偏导不为 0。所以：

$$\frac{\partial J}{\partial z_{[x,y]}^{[l]}} = \sum_{x'=0}^{n_H^{[l+1]}} \sum_{y'=0}^{n_W^{[l+1]}} \frac{\partial J}{\partial z_{[x',y']}^{[l+1]}} g' \left( z_{[x,y]}^{[l]} \right) w_{[a,b]}^{[l+1]} \quad (35)$$

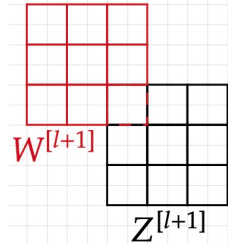
代入  $a = x - x'$ ,  $b = y - y'$  得到：

$$\frac{\partial J}{\partial z_{[x,y]}^{[l]}} = \sum_{x'=0}^{n_H^{[l+1]}} \sum_{y'=0}^{n_W^{[l+1]}} \frac{\partial J}{\partial z_{[x',y']}^{[l+1]}} g' \left( z_{[x,y]}^{[l]} \right) w_{[x-x',y-y']}^{[l+1]} \quad (36)$$

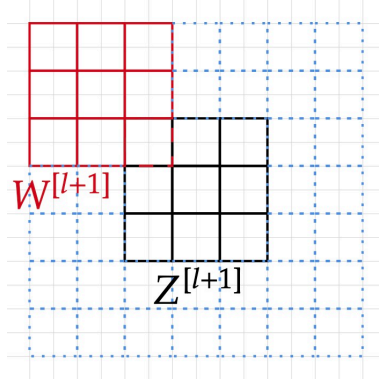
因为是标量乘法所以可以移项得到：

$$\begin{aligned} \frac{\partial J}{\partial z_{[x,y]}^{[l]}} &= \sum_{x'=0}^{n_H^{[l+1]}} \sum_{y'=0}^{n_W^{[l+1]}} \frac{\partial J}{\partial z_{[x',y']}^{[l+1]}} \cdot w_{[x-x',y-y']}^{[l+1]} \cdot g'(z_{[x,y]}^{[l]}) = \frac{\partial J}{\partial \mathbf{Z}^{[l+1]}} \star_{full} \mathbf{W}^{[l+1]} \odot g'(\mathbf{Z}^{[l]}) \\ &= \frac{\partial J}{\partial \mathbf{Z}^{[l+1]}} \star_{full} \text{ROT180}(\mathbf{W}^{[l+1]}) \odot g'(\mathbf{Z}^{[l]}) \end{aligned} \quad (37)$$

注意其中  $\star_{full}$  提示这是一个 full convolution，为什么呢？例如假设  $[x', y'] = [0, 0]$ ,  $[x, y] = [2, 2]$ ，则卷积运算要在  $\mathbf{Z}^{[l+1]}$  左上角覆盖权重矩阵的  $[x - x', y - y'] = [2, 2]$  索引。如图：

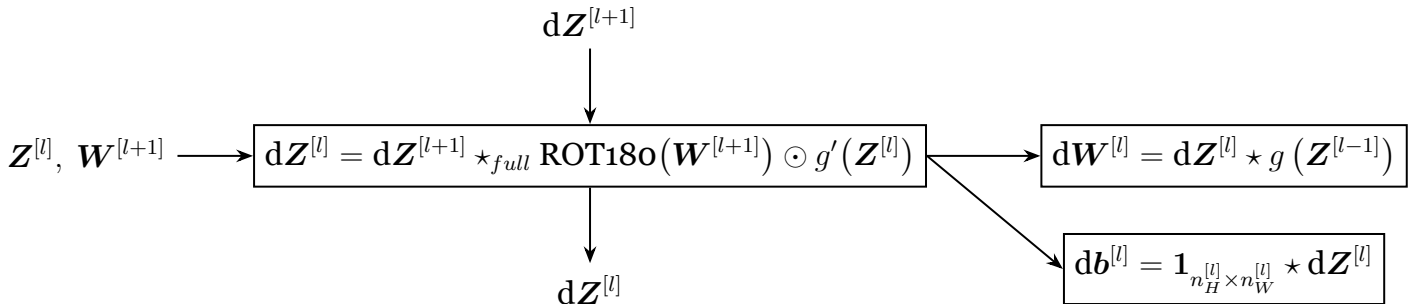


即在卷积运算中  $\mathbf{W}^{[l+1]}$  不能局限在  $\mathbf{Z}^{[l+1]}$  范围中，必须在  $\mathbf{Z}^{[l+1]}$  四周额外用 0 垫衬 (pad)  $n_H^{[l+1]} - n_H^{[l]} = f^{[l+1]} - 1$  的大小：



### 3.4 总结：Convolutional Layer 的 Back Propagation

整合之前的推导，我们可以得到：



## 4 Pooling Layer 的 Back Propagation

对于 max pooling，在 back prop 时将某一池化像素的梯度都回传给上一层中池化范围内的最大值所在像素，其余像素梯度为 0。

对于 average pooling，则把梯度的各个池化区域的值取平均后放在上一层对应的池化位置。

