# Gradient Calculation I

- For convolutional layers, recall we had

$$
\text{vec}(S^{m,i})
$$
$$
= \left( \phi(\text{pad}(Z^{m,i}))^T \otimes \mathcal{I}_{d^{m+1}} \right) \text{vec}(W^m) +
$$
$$
\left( \mathbb{1}_{a^m_{\text{conv}} b^m_{\text{conv}}} \otimes \mathcal{I}_{d^{m+1}} \right) \boldsymbol{b}^m
$$

# Gradient Calculation II

Thus

$$\frac{\partial \xi_i}{\partial \text{vec}(W^m)^T} = \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \frac{\partial \text{vec}(S^{m,i})}{\partial \text{vec}(W^m)^T}$$

$$= \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \left( \phi(\text{pad}(Z^{m,i}))^T \ \otimes \ \mathcal{I}_{d^{m+1}} \right)$$

$$= \text{vec} \left( \frac{\partial \xi_i}{\partial S^{m,i}} \phi(\text{pad}(Z^{m,i}))^T \right)^T \tag{1}$$

# Gradient Calculation III

where (1) is from

$$\text{vec}(AB)^T = \text{vec}(B)^T(\mathcal{I} \otimes A^T) \tag{2}$$
$$= \text{vec}(A)^T(B \otimes \mathcal{I}) \tag{3}$$

- We applied chain rule here
- Note that we define

$$\frac{\partial \mathbf{y}}{\partial(\mathbf{x})^T} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_{|x|}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_{|y|}}{\partial x_1} & \cdots & \frac{\partial y_{|y|}}{\partial x_{|x|}} \end{bmatrix}, \tag{4}$$

# Gradient Calculation IV

where $x$ and $y$ are column vectors, and $|x|$, $|y|$ are their lengths.

- Thus if

$$y = Ax$$

then

$$y_1 = A_{11}x_1 + \cdots + A_{1|x|}x_{|x|}$$

and

$$\frac{\partial y}{\partial (x)^T} = \begin{bmatrix} A_{11} & A_{12} & \cdots \\ A_{21} & & \\ \vdots & & \end{bmatrix} = A$$

# Gradient Calculation V

- Similarly

$$\frac{\partial \xi_i}{\partial (\boldsymbol{b}^m)^T} = \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \frac{\partial \text{vec}(S^{m,i})}{\partial (\boldsymbol{b}^m)^T}$$

$$= \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \left( \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \otimes \mathcal{I}_{d^{m+1}} \right)$$

$$= \text{vec} \left( \frac{\partial \xi_i}{\partial S^{m,i}} \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \right)^T, \qquad (5)$$

where (5) is from (3).

# Gradient Calculation VI

- To calculate (1), $\phi(\text{pad}(Z^{m,i}))$ has been available from the forward process of calculating the function value.

- In (1) and (5), $\partial \xi_i / \partial S^{m,i}$ is also needed

- We will show that it can be obtained by a backward process.

# Calculation of $\partial \xi_i / \partial S^{m,i}$ I

- What we will do is to assume that

$$\frac{\partial \xi_i}{\partial Z^{m+1,i}}$$

  is available

- Then we show details of calculating

$$\frac{\partial \xi_i}{\partial S^{m,i}} \text{ and } \frac{\partial \xi_i}{\partial Z^{m,i}}$$

  for layer $m$.

- Thus a back propagation process

- We have the following workflow.

$$Z^{m,i} \leftarrow \text{padding} \leftarrow \text{convolution} \leftarrow \sigma(S^{m,i})$$
$$\leftarrow \text{pooling} \leftarrow Z^{m+1,i}. \tag{6}$$

- From chain rule,

$$\frac{\partial\xi_i}{\partial\text{vec}(S^{m,i})^T} = \frac{\partial\xi_i}{\partial\text{vec}(\sigma(S^{m,i}))^T}\frac{\partial\text{vec}(\sigma(S^{m,i}))}{\partial\text{vec}(S^{m,i})^T}$$

- If $\sigma$ is a scalar function, then

$$\frac{\partial \text{vec}(\sigma(S^{m,i}))}{\partial \text{vec}(S^{m,i})^T}$$

is a squared diagonal matrix of

$$|\text{vec}(S^{m,i})| \times |\text{vec}(S^{m,i})|$$

# Calculation of $\partial \xi_i / \partial S^{m,i}$ IV

- We further assume that the RELU activation function is used. Recall that we assume

$$\sigma'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

though $\sigma(x)$ is not differentiable at $x = 0$

# Calculation of $\partial \xi_i / \partial S^{m,i}$ V

- We can define

$$I[S^{m,i}]_{(p,q)} = \begin{cases} 1 & \text{if } S^{m,i}_{(p,q)} > 0, \\ 0 & \text{otherwise,} \end{cases}$$

and have

$$\frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} = \frac{\partial \xi_i}{\partial \text{vec}(\sigma(S^{m,i}))^T} \odot \text{vec}(I[S^{m,i}])^T$$

where $\odot$ is Hadamard product (i.e., element-wise products)

- Q: can we extend this to other scalar activation functions?

- Yes, the general form is

$$\frac{\partial\xi_i}{\partial\mathsf{vec}(S^{m,i})^T} = \frac{\partial\xi_i}{\partial\mathsf{vec}(\sigma(S^{m,i}))^T} \odot \mathsf{vec}(\sigma'(S^{m,i}))^T$$

- Next,

$$\frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T}$$

$$= \frac{\partial \xi_i}{\partial \text{vec}(Z^{m+1,i})^T} \frac{\partial \text{vec}(Z^{m+1,i})}{\partial \text{vec}(\sigma(S^{m,i}))^T} \frac{\partial \text{vec}(\sigma(S^{m,i}))}{\partial \text{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial \xi_i}{\partial \text{vec}(Z^{m+1,i})^T} \frac{\partial \text{vec}(Z^{m+1,i})}{\partial \text{vec}(\sigma(S^{m,i}))^T} \right) \odot \text{vec}(I[S^{m,i}])^T$$

$$= \left( \frac{\partial \xi_i}{\partial \text{vec}(Z^{m+1,i})^T} P^{m,i}_{\text{pool}} \right) \odot \text{vec}(I[S^{m,i}])^T \qquad (7)$$

# Calculation of $\partial \xi_i / \partial S^{m,i}$ VIII

- Note that (7) is from

$$Z^{m+1,i} = \mathrm{mat}(P_{\text{pool}}^{m,i} \mathrm{vec}(\sigma(S^{m,i})))_{d^{m+1} \times a^{m+1} b^{m+1}}$$

- If a general scalar activation function is considered, (7) is changed to

$$\frac{\partial \xi_i}{\partial \mathrm{vec}(S^{m,i})^T}$$
$$= \left( \frac{\partial \xi_i}{\partial \mathrm{vec}(Z^{m+1,i})^T} \ P_{\text{pool}}^{m,i} \right) \ \odot \ \mathrm{vec}(\sigma'(S^{m,i}))^T$$

- In the end we calculate $\partial \xi_i / \partial Z^{m,i}$ and pass it to the previous layer.

# Calculation of $\partial \xi_i / \partial S^{m,i}$ X

$$\frac{\partial \xi_i}{\partial \text{vec}(Z^{m,i})^T}$$

$$= \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \frac{\partial \text{vec}(S^{m,i})}{\partial \text{vec}(\phi(\text{pad}(Z^{m,i})))^T} \frac{\partial \text{vec}(\phi(\text{pad}(Z^{m,i})))}{\partial \text{vec}(\text{pad}(Z^{m,i}))^T}$$

$$\frac{\partial \text{vec}(\text{pad}(Z^{m,i}))}{\partial \text{vec}(Z^{m,i})^T}$$

$$= \frac{\partial \xi_i}{\partial \text{vec}(S^{m,i})^T} \left( \mathcal{I}_{a_{\text{conv}}^m b_{\text{conv}}^m} \otimes W^m \right) P_\phi^m P_{\text{pad}}^m \qquad (8)$$

$$= \text{vec} \left( (W^m)^T \frac{\partial \xi_i}{\partial S^{m,i}} \right)^T P_\phi^m P_{\text{pad}}^m, \qquad (9)$$

where (8) is from

$$\text{vec}(S^{m,i})$$
$$= \left( \mathcal{I}_{a^m_{\text{conv}} b^m_{\text{conv}}} \otimes W^m \right) \text{vec}(\phi(\text{pad}(Z^{m,i}))) +$$
$$\left( \mathbb{1}_{a^m_{\text{conv}} b^m_{\text{conv}}} \otimes \mathcal{I}_{d^{m+1}} \right) \boldsymbol{b}^m$$

and (9) is from (2).