

2.

$$\begin{aligned}
 \frac{\partial y_L}{\partial b_i} &= \frac{\partial j_L}{\partial j_{L-1}} \cdot \frac{\partial j_{L-1}}{\partial j_{L-2}} \cdots \frac{\partial j_{i+1}}{\partial j_i} \cdot \frac{\partial j_i}{\partial \tilde{j}_i} \cdot \frac{\partial \tilde{j}_i}{\partial b_i} \\
 &= A_L \prod_{k=i+1}^{L-1} \left[\text{diag}(\sigma'(\tilde{j}_k)) A_k \right] \cdot \text{diag}(\sigma'(y_i)) \\
 &\leq A_L \cdot A_{L-1} \cdots A_{j+1} \cdots A_{i+1} \quad " \leq " \text{ denotes all elements have smaller absolute value than RHS}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial y_L}{\partial A_i} &= \frac{\partial j_L}{\partial j_{L-1}} \cdot \frac{\partial j_{L-1}}{\partial j_{L-2}} \cdots \frac{\partial j_{i+1}}{\partial j_i} \cdot \frac{\partial j_i}{\partial \tilde{j}_i} \cdot \frac{\partial \tilde{j}_i}{\partial A_i} \\
 &= A_L \prod_{k=i+1}^{L-1} \left[\text{diag}(\sigma'(\tilde{j}_k)) A_k \right] \cdot \text{diag}(\sigma'(\tilde{j}_i)) \cdot y_{i-1}^T \\
 &\leq A_L \cdot A_{L-1} \cdots A_{j+1} \cdots A_{i+1} \cdot y_{i-1}^T
 \end{aligned}$$

Since A_j is small and other matrices are not too large,
the two gradients become small.

$$\begin{aligned}
 \frac{\partial y_L}{\partial b_i} &= \frac{\partial j_L}{\partial j_{L-1}} \cdot \frac{\partial j_{L-1}}{\partial j_{L-2}} \cdots \frac{\partial j_{i+1}}{\partial j_i} \cdot \frac{\partial j_i}{\partial \tilde{j}_i} \cdot \frac{\partial \tilde{j}_i}{\partial b_i} \\
 &= A_L \prod_{k=i+1}^{L-1} \left[\text{diag}(\sigma'(\tilde{j}_k)) A_k \right] \cdot \text{diag}(\sigma'(y_i)) \\
 &\leq A_L \cdots A_{j+1} \text{diag}(\sigma'(\tilde{y}_j)) A_j A_{j-1} \cdots A_{i+1}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial y_L}{\partial A_i} &= \frac{\partial j_L}{\partial j_{L-1}} \cdot \frac{\partial j_{L-1}}{\partial j_{L-2}} \cdots \frac{\partial j_{i+1}}{\partial j_i} \cdot \frac{\partial j_i}{\partial \tilde{j}_i} \cdot \frac{\partial \tilde{j}_i}{\partial A_i} \\
 &= A_L \prod_{k=i+1}^{L-1} \left[\text{diag}(\sigma'(\tilde{j}_k)) A_k \right] \cdot \text{diag}(\sigma'(\tilde{j}_i)) \cdot y_{i-1}^T \\
 &\leq A_L \cdot A_{L-1} \cdots A_{j+1} \text{diag}(\sigma'(\tilde{y}_j)) A_j A_{j-1} \cdots A_{i+1} \cdot y_{i-1}^T
 \end{aligned}$$

$$|\tilde{y}_j| \text{ large} \Rightarrow \sigma'(\tilde{y}_j) \text{ small}$$

and the other matrices are not too large.

Hence, the two gradients become small.



3.

$$\begin{aligned}\theta'_I &= \theta^{\circ} - \alpha g^{\circ} + \beta(\theta^{\circ} - \theta^{-1}) \\ &= \theta^{\circ} - \alpha g^{\circ} = \theta^{\circ} - \alpha(g^{\circ} + \beta v^{\circ}) = \theta'_{II}\end{aligned}$$

Suppose that $\theta_I^k = \theta_{II}^k$ for some $k \geq 1$.

That is, $\theta_I^{k+1} = \theta^{k+1} - \alpha g^{k+1} + \beta(\theta^{k+1} - \theta^{k+2})$
 $= \theta^{k+1} - \alpha v^k = \theta_{II}^k$

Then, $\theta^k - \theta^{k+1} = -\alpha v^k$

$$\begin{aligned}\theta_I^{k+1} &= \theta^k - \alpha g^k + \beta(\theta^k - \theta^{k+1}) \\ &= \theta^k - \alpha g^k + \beta(-\alpha v^k) \\ &= \theta^k - \alpha(g^k + \beta v^k) = \theta_{II}^{k+1}\end{aligned}$$

∴ By induction, form I and II are identical.

□

4.

i) $y_1[k, i, j]$ depends on $X[:, i-2:i+2, j-2:j+2]$.
 $i, j \leq 0$ or $i, j \geq 225$ is due to zero padding
and can be ignored.

ii) $y_2[k, i, j]$ depends on $y_1[:, 2i-1:2i, 2j-1:2j]$,
so on $X[:, 2i-3:2i+2, 2j-3:2j+2]$ by i).

iii) $y_3[k, i, j]$ depends on $y_2[:, 2i-3:2i+2, 2j-3:2j+2]$ by ii),
so on $X[:, 4i-9:4i+6, 4j-9:4j+6]$ by ii).

□

5.

(i) First module has

$$\begin{aligned}
 & ((1 \times 1) \times 256 + 1) \times 128 \\
 & + ((3 \times 3) \times 256 + 1) \times 128 \\
 & + ((5 \times 5) \times 256 + 1) \times 128 = 1,147,264 \text{ trainable parameters.}
 \end{aligned}$$

Second module has

$$\begin{aligned}
 & ((1 \times 1) \times 256 + 1) \times 128 \\
 & + ((1 \times 1) \times 256 + 1) \times 64 + ((3 \times 3) \times 64 + 1) \times 192 \\
 & + ((1 \times 1) \times 256 + 1) \times 64 + ((5 \times 5) \times 64 + 1) \times 96 \\
 & + ((1 \times 1) \times 256 + 1) \times 64 = 346,720 \text{ trainable parameters.}
 \end{aligned}$$

The second module has much less trainable parameters.

(ii)

	module 1	module 2
addition	$ \begin{aligned} & 256 \times 32 \times 32 \times 128 \\ & + 256 \times 3 \times 3 \times 32 \times 32 \times 192 \\ & + 256 \times 5 \times 5 \times 32 \times 32 \times 96 \\ & = 1,115,684,864 \end{aligned} $	$ \begin{aligned} & 256 \times 32 \times 32 \times 128 \\ & + 256 \times 32 \times 32 \times 64 + 64 \times 3 \times 3 \times 32 \times 32 \times 192 \\ & + 256 \times 32 \times 32 \times 64 + 64 \times 5 \times 5 \times 32 \times 32 \times 96 \\ & + 256 \times 32 \times 32 \times 64 \\ & = 354,418,688 \end{aligned} $
multiplication	$ \begin{aligned} & 256 \times 32 \times 32 \times 128 \\ & + 256 \times 3 \times 3 \times 32 \times 32 \times 192 \\ & + 256 \times 5 \times 5 \times 32 \times 32 \times 96 \\ & = 1,115,684,864 \end{aligned} $	$ \begin{aligned} & 256 \times 32 \times 32 \times 128 \\ & + 256 \times 32 \times 32 \times 64 + 64 \times 3 \times 3 \times 32 \times 32 \times 192 \\ & + 256 \times 32 \times 32 \times 64 + 64 \times 5 \times 5 \times 32 \times 32 \times 96 \\ & + 256 \times 32 \times 32 \times 64 \\ & = 354,418,688 \end{aligned} $
activation function evaluation	$ \begin{aligned} & 32 \times 32 \times 128 \\ & + 32 \times 32 \times 192 \\ & + 32 \times 32 \times 96 \\ & = 425,984 \end{aligned} $	$ \begin{aligned} & 32 \times 32 \times 128 \\ & + 32 \times 32 \times 64 + 32 \times 32 \times 192 \\ & + 32 \times 32 \times 64 + 32 \times 32 \times 96 \\ & + 32 \times 32 \times 64 \\ & = 622,592 \end{aligned} $

There are much more add, multiplying in module 1.

Activation function evaluation is slightly more in module 2.

□