

1. $w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \in \mathbb{R}^{2 \times 3 \times 3}$, $w_1, w_2 \in \mathbb{R}^{3 \times 3}$

$$w_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad w_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

□

2. AvgPool2d can be represented as a convolution

with a filter $w = \begin{bmatrix} w_1 \\ \vdots \\ w_c \end{bmatrix} \in \mathbb{R}^{C \times C \times K \times K}$

$$w_i = [\underbrace{0, \dots, 0}_\text{KxK}, \underbrace{\begin{bmatrix} \frac{1}{K^2} & \dots & \frac{1}{K^2} \\ \vdots & \ddots & \vdots \\ \frac{1}{K^2} & \dots & \frac{1}{K^2} \end{bmatrix}, 0, \dots, 0}_\text{i-th element of } w_i]^T \in \mathbb{R}^{C \times K \times K}, \quad i=1, \dots, C$$

with stride K without zero paddings.

□

3. $w = \begin{bmatrix} [[0.299]], \\ [[0.587]], \\ [[0.114]] \end{bmatrix} \in \mathbb{R}^{3 \times 1 \times 1}$ be the filter.

4. Since $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ nondecreasing,

$$x_1 \geq x_2 \Rightarrow \sigma(x_1) \geq \sigma(x_2), \quad \forall x_1, x_2 \in \mathbb{R}$$

$$\therefore \rho(x) = x_{ij} \Rightarrow \rho(\sigma(x)) = \sigma(x_{ij}), \quad x \in \mathbb{R}^{m \times n}, \quad \text{for some } i, j$$

where p is the kernel size of max pool ρ .

$$\text{That is, } \sigma(\rho(x)) = \sigma(x_{ij}) = \rho(\sigma(x)).$$

Since stride of max pooling ρ is also p ,

$$\text{the equation } \sigma(\rho(x)) = \rho(\sigma(x)) \text{ holds for } \forall x \in \mathbb{R}^{m \times n}$$

□

6.

(a)

$$\frac{\partial y_L}{\partial b_L} = \frac{\partial}{\partial b_L} (A_L y_{L-1} + b_L) = \frac{\partial}{\partial b_L} \begin{bmatrix} (b_L)_1 \\ \vdots \\ (b_L)_{n_L} \end{bmatrix} = I$$

$$\frac{\partial y_L}{\partial y_{L-1}} = \frac{\partial}{\partial y_{L-1}} (A_L y_{L-1} + b_L) = \frac{\partial}{\partial y_{L-1}} [(A_L)_{1,1}(y_{L-1})_1 + \dots + (A_L)_{1,n_{L-1}}(y_{L-1})_{n_{L-1}}]$$

$$= \begin{bmatrix} (A_L)_{1,1} & \dots & (A_L)_{1,n_{L-1}} \end{bmatrix} = A_L$$

$$\begin{aligned} \frac{\partial y_L}{\partial b_L} &= \frac{\partial}{\partial b_L} \sigma(A_L y_{L-1} + b_L) = \frac{\partial}{\partial b_L} \begin{bmatrix} \sigma((A_L)_{1,1}(y_{L-1})_1 + \dots + (A_L)_{1,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) \\ \vdots \\ \sigma((A_L)_{n_L,1}(y_{L-1})_1 + \dots + (A_L)_{n_L,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) \end{bmatrix} \\ &= \begin{bmatrix} \sigma'((A_L)_{1,1}(y_{L-1})_1 + \dots + (A_L)_{1,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) & 0 & \dots & 0 \\ \vdots & & & \\ 0 & \dots & \sigma'((A_L)_{n_L,1}(y_{L-1})_1 + \dots + (A_L)_{n_L,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) & \end{bmatrix} \\ &= \text{diag}(\sigma'(A_L y_{L-1} + b_L)), \quad l = 1, \dots, L-1 \end{aligned}$$

$$\begin{aligned} \frac{\partial y_L}{\partial y_{L-1}} &= \frac{\partial}{\partial y_{L-1}} \sigma(A_L y_{L-1} + b_L) = \frac{\partial}{\partial y_{L-1}} \begin{bmatrix} \sigma((A_L)_{1,1}(y_{L-1})_1 + \dots + (A_L)_{1,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) \\ \vdots \\ \sigma((A_L)_{n_L,1}(y_{L-1})_1 + \dots + (A_L)_{n_L,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) \end{bmatrix} \\ &= \begin{bmatrix} \sigma'((A_L)_{1,1}(y_{L-1})_1 + \dots + (A_L)_{1,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) & 0 & \dots & 0 \\ \vdots & & & \\ 0 & \dots & \sigma'((A_L)_{n_L,1}(y_{L-1})_1 + \dots + (A_L)_{n_L,n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) & \end{bmatrix} \begin{bmatrix} (A_L)_{1,1} & \dots & (A_L)_{1,n_{L-1}} \\ \vdots & & \vdots \\ (A_L)_{n_L,1} & \dots & (A_L)_{n_L,n_{L-1}} \end{bmatrix} \\ &= \text{diag}(\sigma'(A_L y_{L-1} + b_L)) A_L, \quad l = 2, \dots, L-1 \end{aligned}$$

$$\begin{aligned}
 (b) \quad \left(\frac{\partial y_L}{\partial A_L} \right)_{ij} &= \frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial}{\partial (A_L)_{ij}} (A_L y_{L-1} + b_L) \\
 &= \frac{\partial}{\partial (A_L)_{ij}} \left((A_L)_{11} (y_{L-1})_1 + \dots + (A_L)_{n_{L-1}} (y_{L-1})_{n_{L-1}} \right) \\
 &= (y_{L-1})_j \\
 \therefore \frac{\partial y_L}{\partial A_L} &= J_{L-1}^T
 \end{aligned}$$

$$\begin{aligned}
 \left(\frac{\partial y_L}{\partial A_L} \right)_{ij} &= \frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial y_L}{\partial y_L} \frac{\partial y_L}{\partial (A_L)_{ij}} \\
 &= \frac{\partial y_L}{\partial y_L} \frac{\partial}{\partial (A_L)_{ij}} \sigma(A_L y_{L-1} + b_L) \\
 &= \frac{\partial y_L}{\partial y_L} \frac{\partial}{\partial (A_L)_{ij}} \left[\begin{array}{c} \sigma((A_L)_{11}(y_{L-1})_1 + \dots + (A_L)_{1n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) \\ \vdots \\ \sigma((A_L)_{n_{L-1}1}(y_{L-1})_1 + \dots + (A_L)_{n_L n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) \end{array} \right] \\
 &= \frac{\partial y_L}{\partial y_L} \left[\begin{array}{c} \circ \\ \sigma'((A_L)_{11}(y_{L-1})_1 + \dots + (A_L)_{1n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) \cdot (y_{L-1})_1 \\ \vdots \\ \circ \\ \sigma'((A_L)_{n_{L-1}1}(y_{L-1})_1 + \dots + (A_L)_{n_L n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_{n_L}) \cdot (y_{L-1})_{n_L} \end{array} \right] \\
 &= \sigma'((A_L)_{11}(y_{L-1})_1 + \dots + (A_L)_{1n_{L-1}}(y_{L-1})_{n_{L-1}} + (b_L)_1) \left(\frac{\partial y_L}{\partial y_L} \right)_i (y_{L-1})_i
 \end{aligned}$$

$$\therefore \frac{\partial y_L}{\partial A_L} = \text{diag}(\sigma'(A_L y_{L-1} + b_L)) \left(\frac{\partial y_L}{\partial y_L} \right)^T J_{L-1}^T$$

□