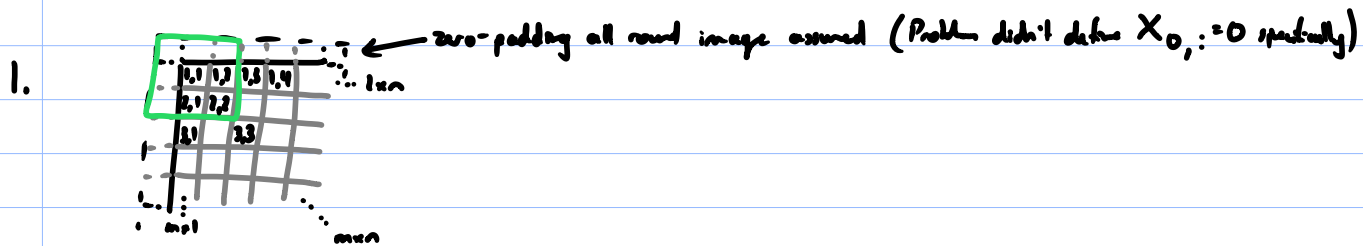


MFDNN

Homework 4

24/03/26



Deriv:  $\sum \begin{pmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{pmatrix} \odot \begin{pmatrix} 0 & 0 & 0 \\ 0 & x_{11} & x_{12} \\ 0 & x_{21} & x_{22} \end{pmatrix} = X_{21} - X_{11}$  Suggests  $w_4 = 1, w_5 = -1, w_i = 0 \forall i \neq 4, 5$

$\sum \begin{pmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{pmatrix} \odot \begin{pmatrix} 0 & 0 & 0 \\ 0 & x_{11} & x_{12} \\ 0 & x_{21} & x_{22} \end{pmatrix} = X_{12} - X_{11}$   $w_6 = 1, w_5 = -1, w_i = 0 \forall i \neq 5, 6$

So for  $w \in \mathbb{R}^{2 \times 3 \times 3}$ :

$$w = \left( \begin{pmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{pmatrix} \right)$$

Sanity check:

$$Y_{l,i,j} = \sum_{a=1}^3 \sum_{b=1}^3 w_{l,a,b} X'_{i+a-1,j+b-1}$$

← padded X

$$Y_{1,1,1} = w_{22} X_{11} + w_{32} X_{21} = -X_{11} + X_{21} \checkmark$$

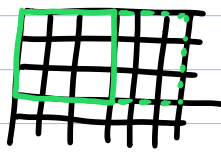
$$Y_{2,1,1} = w_{22} X_{11} + w_{32} X_{12} = -X_{11} + X_{12} \checkmark$$

$$Y_{1,2,2} = w_{22} X_{22} + w_{32} X_{32} = -X_{22} + X_{32} \checkmark$$

2. The Avg Pool2d operation conceptually simply outputs the average of the 'window' being considered.

Since a convolution outputs the sum of the product of the kernel weights with the window elements simply define  $w = \frac{1}{k^2} \mathbf{1} \in \mathbb{R}^{k \times k}$  and use stride of  $k$  to ensure no overlaps of kernel.

↳ or simply take  $w$  as a constant multiplier/kernel



3.  $w = ((0.299), (0.587), (0.114))$

↑  
 $\mathbb{R}^{1 \times 1}$

4. Consider a block,  $B \in \mathbb{R}^{p \times q}$  ( $p \leq m, q \leq n$ ) in  $X \in \mathbb{R}^{m \times n}$  where  $p$  and  $q$  are the kernel size of the max pool. Let  $X_{\max}$  be <sup>(one of)</sup> the largest element of  $B$ .  $\rho(B) = X_{\max}$  by definition of  $X_{\max}$ . Now we apply the activation function to this block's output to give,  $\sigma(\rho(B)) = \sigma(X_{\max})$ .

Now,  $X_{\max} \geq X_{ij}$  for any  $X_{ij} \in B$  by definition of  $X_{\max}$ .

Since  $\sigma$  is non-decreasing,  $X_{\max} \geq X_{ij} \Rightarrow \sigma(X_{\max}) \geq \sigma(X_{ij})$  for all  $X_{ij} \in B$ .

So if we were to apply the max pool operation now on  $\sigma(B)$ , clearly,  $\rho(\sigma(B)) = \sigma(X_{\max})$ . since  $\sigma(X_{\max})$  is also the largest element of  $\sigma(B)$  (by above).

Since this argument is true for any and all blocks  $B \in X$ , this implies  $\rho(\sigma(X)) = \sigma(\rho(X))$ .  
 $\hookrightarrow$  and hence all elements of output  $\mathbb{R}^{k \times l}$  matrix

5. (Workings)

$$f_i(\underline{a}, \underline{b}) := \|\underline{P}(y_i) - \underline{f}_{\underline{a}, \underline{b}}(\underline{x}_i)\|^2$$

$$\underline{z} := \underline{a}^T \underline{x} + b$$

$$= \|\underline{P}(y_i) - \{\sigma(-\underline{z}), \sigma(\underline{z})\}\|^2$$

$$= \underbrace{(1-y_i) \frac{1}{2}}_{\substack{\text{if } y_i = -1, = 1 \\ \text{if } y_i = 1, = 0}} \|\{1, 0\} - \{\sigma(-\underline{z}), \sigma(\underline{z})\}\|^2 + \underbrace{(1+y_i) \frac{1}{2}}_{\substack{\text{if } y_i = 1, = 1 \\ \text{if } y_i = -1, = 0}} \|\{0, 1\} - \{\sigma(-\underline{z}), \sigma(\underline{z})\}\|^2$$

$$= \frac{1}{2}(1-y_i) \left( (1-\sigma(-\underline{z}))^2 + \sigma(\underline{z})^2 \right) + \frac{1}{2}(1+y_i) \left( \sigma(-\underline{z})^2 + (1-\sigma(\underline{z}))^2 \right)$$

$$=: \ell(\underline{z}, y_i) = \ell(\underline{a}^T \underline{x} + b, y_i)$$

$$6.(a) \quad \frac{\partial y_L}{\partial b_L} = \frac{\partial}{\partial b_L} (A_L y_{L-1} + b_L) = 1$$

$$\frac{\partial y_L}{\partial y_{L-1}} = \frac{\partial}{\partial y_{L-1}} (A_L y_{L-1} + b_L) = A_L$$

For  $l=1, \dots, L-1$ ,

$$\frac{\partial y_l}{\partial b_l} = \begin{pmatrix} \frac{\partial y_{l1}}{\partial b_{l1}} & \frac{\partial y_{l1}}{\partial b_{l2}} & \dots & \frac{\partial y_{l1}}{\partial b_{ln_l}} \\ \frac{\partial y_{l2}}{\partial b_{l1}} & \frac{\partial y_{l2}}{\partial b_{l2}} & \dots & \frac{\partial y_{l2}}{\partial b_{ln_l}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{ln_l}}{\partial b_{l1}} & \frac{\partial y_{ln_l}}{\partial b_{l2}} & \dots & \frac{\partial y_{ln_l}}{\partial b_{ln_l}} \end{pmatrix}$$

Define  $\tilde{y}_l = A_l y_{l-1} + b_l$  (i.e.  $y_l$  pre-activation)

$$\begin{aligned} \frac{\partial \tilde{y}_{li}}{\partial b_{lj}} &= \frac{\partial}{\partial b_{lj}} ([A_l y_{l-1} + b_l]_i) \\ &= \frac{\partial}{\partial b_{lj}} ([A_l y_{l-1}]_i + b_{li}) \\ &= \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} \end{aligned}$$

Hence, only diagonal entries of  $\frac{\partial y_l}{\partial b_l}$  remain (i.e.  $i=j$ ) are non-zero

$$\begin{aligned} \frac{\partial y_{li}}{\partial b_{lj}} &= \frac{\partial}{\partial b_{lj}} ([\sigma(\tilde{y}_l)]_i) \\ &= \frac{\partial \tilde{y}_{li}}{\partial b_{lj}} [\sigma'(\tilde{y}_l)]_i = \frac{\partial \tilde{y}_{li}}{\partial b_{lj}} \sigma'(\tilde{y}_{li}) \end{aligned}$$

by chain rule and elementwise operation

$$\text{with } \frac{\partial y_{li}}{\partial b_{li}} = \sigma'(\tilde{y}_{li}) = \sigma'([A_l y_{l-1} + b_l]_i)$$

$$= \begin{cases} 0 & \text{if } i \neq j \\ \sigma'(\tilde{y}_{li}) & \text{if } i = j \end{cases} \text{ by above}$$

$$\text{That is, } \frac{\partial y_l}{\partial b_l} = \text{diag}(\sigma'(A_l y_{l-1} + b_l)) \quad \square$$

For  $l=2, \dots, L-1$ ,

For  $\frac{\partial y_l}{\partial y_{l-1}}$ , consider entries one row:

$$\begin{aligned} \frac{\partial \tilde{y}_{li}}{\partial y_{l-1j}} &= \frac{\partial}{\partial y_{l-1j}} ([A_l y_{l-1}]_i + b_{li}) \\ &= \frac{\partial}{\partial y_{l-1j}} \left( \sum_{k=1}^{n_{l-1}} (A_l)_{ik} (y_{l-1})_k + b_{li} \right) \\ &= (A_l)_{ij} \end{aligned}$$

$$i \left( \overbrace{\quad}^{A_l} \right) \left( \underbrace{\quad}_{y_{l-1}} \right)$$

$$\begin{aligned} \text{So } \frac{\partial y_{li}}{\partial y_{l-1j}} &= \frac{\partial}{\partial y_{l-1j}} (\sigma(\tilde{y}_{li})) = \frac{\partial \tilde{y}_{li}}{\partial y_{l-1j}} \sigma'(\tilde{y}_{li}) = (A_l)_{ij} \sigma'(\tilde{y}_{li}) \\ &= \sigma'(\tilde{y}_{li}) (A_l)_{ij} \\ &= [\sigma'(A_l y_{l-1} + b_l)]_i (A_l)_{ij} \end{aligned}$$

elements equivalent

$$\begin{aligned} [\text{diag}(\sigma'(A_l y_{l-1} + b_l)) A_l]_{ij} &= \sum_{k=1}^{n_l} [\text{diag}(\sigma'(A_l y_{l-1} + b_l))]_{ik} (A_l)_{kj} = [\sigma'(A_l y_{l-1} + b_l)]_i (A_l)_{ij} \\ &\quad \underbrace{\neq 0 \text{ iff } k=i} \\ \Rightarrow \frac{\partial y_l}{\partial y_{l-1}} &= \text{diag}(\sigma'(A_l y_{l-1} + b_l)) A_l \end{aligned}$$

$$i \left( \overbrace{\quad}^{A_l} \right) \left( \underbrace{\quad}_{y_{l-1}} \right)$$

as given  $\square$

$$6.(b) \left( \frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial}{\partial (A_L)_{ij}} (A_L y_{L-1} + b_L) \quad A_L \in \mathbb{R}^{n_L \times n_{L-1}}$$

so only  $i=1$  defined

$$= \frac{\partial}{\partial (A_L)_{ij}} \left( \sum_{k=1}^{n_{L-1}} (A_L)_{i,k} (y_{L-1})_k \right)$$

$$= (y_{L-1})_j \Rightarrow \frac{\partial y_L}{\partial A_L} = ((y_{L-1})_1, (y_{L-1})_2, \dots, (y_{L-1})_{n_{L-1}})$$

$$= y_{L-1}^T$$

For  $l=1, \dots, L-1$ ,

$$\frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial y_L}{\partial y_L} \frac{\partial y_L}{\partial (A_L)_{ij}} = \left( \frac{\partial y_L}{\partial y_{L-1}} \frac{\partial y_{L-1}}{\partial y_{L-2}} \dots \frac{\partial y_{L-1}}{\partial y_L} \right) \frac{\partial y_L}{\partial (A_L)_{ij}} \quad \text{by chain rule}$$

$$= (A_L) (\text{diag}(\sigma'(A_L y_{L-1} + b_L)) A_L) (\text{diag}(\sigma'(A_{L-1} y_{L-2} + b_{L-1})) A_{L-1}) \dots (\text{diag}(\sigma'(A_{L+1} y_L + b_{L+1})) A_{L+1}) \frac{\partial y_L}{\partial (A_L)_{ij}} \quad \text{by part (a)}$$

$$\frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial}{\partial (A_L)_{ij}} \left( \sigma \left[ \sum_{p=1}^{n_L} \left( s_p + \sum_{k=1}^{n_{L-1}} (A_L)_{p,k} (y_{L-1})_k \right) + b_L \right] \right)$$

$$= s_i (y_{L-1})_j \sigma'(\tilde{y}_L)$$













