

MFDNN

Homework 7

24/04/26

$$1.(a) \lim_{\beta \rightarrow \infty} \left(\frac{1}{\beta} \log \left(\sum_{i=1}^n \exp(\beta x_i) \right) \right) = \lim_{\beta \rightarrow \infty} \left(\frac{1}{\beta} \left(\log(\exp(\beta x_j)) + \log \left(1 + \sum_{\substack{i=1, \\ i \neq j}}^n \frac{\exp(\beta x_i)}{\exp(\beta x_j)} \right) \right) \right)$$

$$\text{since } \log(a+b) = \log(a) + \log\left(1 + \frac{b}{a}\right); \text{ where } x_j := \max\{x_1, \dots, x_n\}$$

$$= \lim_{\beta \rightarrow \infty} \left(x_j + \frac{1}{\beta} \log \left(1 + \sum_{\substack{i=1, \\ i \neq j}}^n \exp(\beta(x_i - x_j)) \right) \right)$$

Since x_j is the max x , $x_i - x_j \leq 0 \forall i$

\Rightarrow argument of exponential is negative

So in limit $\beta \rightarrow \infty$, $\exp(\beta(x_i - x_j)) \rightarrow 0 \forall i$

(if the maximum is not unique, then the summation is a finite integer $< n$)
(each shared max adds 1)

In any case,

$$= \lim_{\beta \rightarrow \infty} \left(x_j + \frac{\log(1+\gamma)}{\beta} \right), \gamma \rightarrow 0 \text{ if max unique, else } 0 < \gamma < n$$

$$= x_j + 0$$

$$= \max\{x_1, \dots, x_n\} \text{ as required } \blacksquare$$

$$(b) \nabla \gamma_i = \left(\frac{\partial}{\partial x_1} \gamma_i, \frac{\partial}{\partial x_2} \gamma_i, \dots, \frac{\partial}{\partial x_n} \gamma_i \right) \quad \frac{\partial}{\partial x_j} \gamma_i = \frac{\partial}{\partial x_j} \left(\log \left(\sum_{i=1}^n e^{x_i} \right) \right)$$

$$= \left(\frac{e^{x_1}}{\sum_{i=1}^n e^{x_i}}, \frac{e^{x_2}}{\sum_{i=1}^n e^{x_i}}, \dots, \frac{e^{x_n}}{\sum_{i=1}^n e^{x_i}} \right) \quad = e^{x_j} \cdot \frac{1}{\sum_{i=1}^n e^{x_i}}$$

$$\text{Clearly, } (\nabla \gamma_i(\underline{x}))_k = \frac{\exp(x_k)}{\sum_{i=1}^n \exp(x_i)} = \mu_k(\underline{x})$$

$$\text{True } \forall k \Rightarrow \nabla \gamma_i = \mu \blacksquare$$

$$(c) \frac{\partial}{\partial x_j} \gamma_\beta = \frac{1}{\beta} \cdot \frac{\beta \exp(\beta x_j)}{\sum_{i=1}^n \exp(\beta x_i)}$$

$$\nabla \gamma_\beta(\underline{x}) = \sum_{j=1}^n \frac{\exp(\beta x_j)}{\sum_{i=1}^n \exp(\beta x_i)} \underline{e}_j = \sum_{j=1}^n \frac{1}{\sum_{i=1}^n \exp(\beta(x_i - x_j))} \underline{e}_j$$

• For term where $j = i_{\max}$, $x_i - x_j < 0 \Rightarrow \exp(\beta(x_i - x_j)) \rightarrow 0$ as $\beta \rightarrow \infty$
for all $i \neq j$

$$\text{For } i=j, x_i - x_j = 0 \Rightarrow \exp(\beta(x_i - x_j)) = 1$$

So for this term, denominator is 1 and so $\underline{e}_j = \underline{e}_{i_{\max}}$ term is kept.

• For all other j , $\exists i \text{ s.t. } x_i > x_j \Rightarrow x_i - x_j > 0 \Rightarrow \exists i \text{ for which } \exp(\beta(x_i - x_j)) \rightarrow \infty$
as $\beta \rightarrow \infty$

So the denominator tends to infinity and the term vanishes and \underline{e}_j is lost.

$$\text{So as } \beta \rightarrow \infty, \nabla \gamma_\beta(\underline{x}) \rightarrow \underline{e}_{i_{\max}} \blacksquare \text{ (i.e. only } j = i_{\max} \text{ term } (\underline{e}_{i_{\max}}) \text{ is preserved)}$$

including additional operations due to padding

2. Num convolution additions = Num convolution multiplications

$$= h_{out}^2 \times k_h^2 \times C_{in} \times C_{out}$$

(HW 5.5)

$$h_{out} = \text{output image width/height} = \left\lfloor \frac{h_{in} - k_h + 2P}{S} + 1 \right\rfloor$$

applicable to padding too

$$\text{Num linear layer multiplications} = C_{out} C_{in} \xrightarrow{\text{summation}} \text{bias}$$

$$\text{Num linear layer additions} = C_{out} (C_{in} - 1 + 1) = C_{out} C_{in}$$

$$y_i = (Ax)_i + b_i \quad i=1, \dots, C_{out} \\ = \sum_{j=1}^{C_{in}} A_{ij} x_j + b_i$$

So we have:

convolution num	C_{in}	C_{out}	k_h	h_{in}	P	S	h_{out}
1	3	64	11	227	0	4	55
2	64	192	5	27	2	1	27
3	192	384	3	13	1	1	13
4	384	256	3	13	1	1	13
5	256	256	3	13	1	1	13

← no padding

linear num	C_{in}	C_{out}
1	9216	4096
2	4096	4096
3	4096	1000

← num classes

← add & mult

$$\Rightarrow \text{Convolutional layers total} = 2 \times (\text{conv1} + \text{conv2} + \text{conv3} + \text{conv4} + \text{conv5})$$

$$= 131133056$$

each $C_{in} \times C_{out} \times k_h^2 \times h_{out}^2$ from table

$$\text{Linear layers total} = 2 \times (\text{lin1} + \text{lin2} + \text{lin3})$$

$$= 117243904$$

3. (Working)

$$\omega \in \mathbb{R}^{\text{normalised } C_{out} \times C_{in} \times f_1 \times f_2}$$

$$Y_{k,\ell,i,j} = \sum_{\gamma=1}^{C_{in}} \sum_{\alpha=1}^{f_1} \sum_{\beta=1}^{f_2} \omega_{\ell,\gamma,\alpha,\beta} X_{k,\gamma,\ell+\alpha-2,j+\beta-2} + b_{\ell}$$

$$\text{out channel } \downarrow \\ BN_{\gamma,\beta}(Y)[b, :, i, j] = \gamma[:, :] \frac{Y[b, :, i, j] - \hat{\mu}[:, :]}{\sqrt{\hat{\sigma}^2[:, :] + \epsilon}} + \beta[:, :]$$

$$= \frac{\gamma[:, :] Y[b, :, i, j]}{\sqrt{\hat{\sigma}^2[:, :] + \epsilon}} - \underbrace{\frac{\gamma[:, :] \hat{\mu}[:, :]}{\sqrt{\hat{\sigma}^2[:, :] + \epsilon}}}_{\text{new bias}} + \beta[:, :]$$

$$Y'_{k,\ell,i,j} = \frac{\gamma[\ell]}{\sqrt{\hat{\sigma}^2[\ell] + \epsilon}} \left(\sum \sum \sum \omega X + b_{\ell} \right)$$

$$= \underbrace{\frac{\gamma[\ell]}{\sqrt{\hat{\sigma}^2[\ell] + \epsilon}}}_{\text{weight multiplier}} \sum \sum \sum \omega X + \underbrace{\frac{\gamma[\ell]}{\sqrt{\hat{\sigma}^2[\ell] + \epsilon}} b_{\ell}}_{\text{to be inc. in bias}}$$

$$A \in \mathbb{R}^{\text{normalised } C_{out} \times C_{in}}$$

precision? question

$$4.(a) \quad \frac{\partial y_L}{\partial y_{L-1}} = \frac{\partial}{\partial y_{L-1}} (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})$$

$$= A_{w_L}$$

$$\left(\frac{\partial y_L}{\partial y_{L-1}} \right)_{ij} = \frac{\partial}{\partial (y_{L-1})_j} \left(\sigma \left(\sum_k (A_{w_L})_{ik} (y_{L-1})_k + b_L \right) \right)$$

$$= (A_{w_L})_{ij} \cdot \sigma'_i (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})$$

$$\Rightarrow \frac{\partial y_L}{\partial y_{L-1}} = \text{diag}(\sigma'(A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})) A_{w_L} \quad \text{(see HW4.6(a) equivalent weights)}$$

$$\left(\frac{\partial y_L}{\partial w_L} \right)_{l,i} = \sum_{k=1}^{n_L} \left(\frac{\partial y_L}{\partial y_L} \right)_{l,k} \left(\frac{\partial y_L}{\partial w_L} \right)_{k,i} \quad i=1, \dots, f_L$$

$$= \sum_{k=1}^{n_L} \left(\frac{\partial y_L}{\partial y_L} \right)_{l,k} \sigma'_k (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L}) (y_{L-1})_{k,i-1}$$

$$= \sum_{k=1}^{n_{L-1}-f_{L-1}+1} \left(\frac{\partial y_L}{\partial y_L} \right)_{l,k} \sigma'_k (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L}) (y_{L-1})_{k,i-1}$$

... ?

$$= \sum_{m=1}^{n_{L-1}} (C_{v_L^T})_{i,m} (y_{L-1})_m = (C_{v_L^T} y_{L-1})_i$$

$$= (C_{v_L^T} y_{L-1})_{l,i}^T$$

$$\frac{\partial y_L}{\partial y_L} \in \mathbb{R}^{1 \times n_L}$$

$$\text{diag}(\sigma'(\dots)) \in \mathbb{R}^{n_L \times n_L}$$

$$y_{L-1} \in \mathbb{R}^{n_{L-1} \times 1}$$

$$v_L \in \mathbb{R}^{1 \times n_L}$$

$$v_L^T \in \mathbb{R}^{n_L \times 1}$$

$$C_{v_L^T} \in \mathbb{R}^{f_L \times n_{L-1}}$$

$$\left(\frac{\partial y_L}{\partial w_L} \right)_{ij} = \frac{\partial (y_L)_i}{\partial (w_L)_j} = \frac{\partial}{\partial (w_L)_j} \left(\sigma \left(\sum_{k=1}^{n_{L-1}} (A_{w_L})_{ik} (y_{L-1})_k + b_L \right) \right)$$

$$\frac{\partial y_L}{\partial w_L} \in \mathbb{R}^{1 \times f_L} \supset (C_{v_L^T} y_{L-1})^T$$

$$C_{v_L^T} y_{L-1} \in \mathbb{R}^{f_L \times 1}$$

$$= \delta_{k,i+j-1} (y_{L-1})_k \sigma'_i (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})$$

$$= (y_{L-1})_{i+j-1} \sigma'_i (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})$$

$$\frac{\partial (A_{w_L})_{ik}}{\partial (w_L)_j} = \delta_{k,i+j-1}$$

$(w_L)_j$ located in positions in A_{w_L} \leftarrow

$$A_{w_L} = \begin{pmatrix} (w_L)_1 & \dots & (w_L)_{f_L} & 0 & 0 & \dots & 0 \\ 0 & (w_L)_1 & \dots & (w_L)_{f_L} & 0 & \dots & 0 \\ 0 & & \ddots & & \ddots & & \vdots \\ \vdots & & & 0 & (w_L)_1 & \dots & (w_L)_{f_L} \end{pmatrix}$$

$(1,j)$
 $(2,j+1)$
 \vdots
 $(i,j+(i-1))$

(as in HW1.5)
(or Chp3 Slk 19)

(a) (cont.) $\frac{\partial y_L}{\partial b_L} = \frac{\partial y_L}{\partial y_L} \frac{\partial y_L}{\partial b_L}$

$$\left(\frac{\partial y_L}{\partial b_L}\right)_i = \frac{\partial}{\partial b_L} \left(\sigma \left(\sum_k (A_{w_L})_{ik} (y_{L-1})_k + b_L \right) \right)$$

$$= \sigma'_i (A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})$$

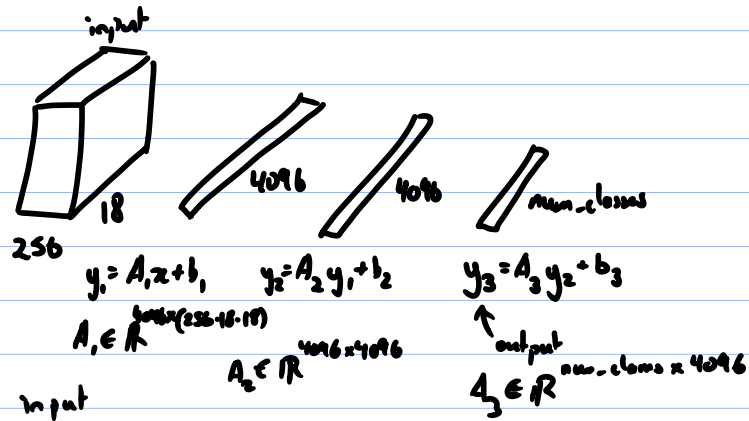
$$= \frac{\partial y_L}{\partial y_L} \text{diag}(\sigma'(A_{w_L} y_{L-1} + b_L \mathbf{1}_{n_L})) \mathbf{1}_{n_L}$$

$$= \underline{v_L} \mathbf{1}_{n_L} \mathbf{0}$$

- (b) In the forward pass, matrix-vector products with A_{w_i} are used to perform convolutions (for $A_{w_i} y_{i-1}$)
- In backpropagation, calculations will require updating $\frac{\partial y_L}{\partial y_{L-1}}$ using right multiplication by A_{w_L} , which should use $A_{w_L}^T$ in transpose-convolutions. (i.e. matrix)
- Calculations will also require $\frac{\partial y_L}{\partial w_L}$ which uses the convolutional operator $C_{v_L^T}$. This should be performed using convolution since part of a matrix-vector product. Note that computing $C_{v_L^T}$ (by extension v_L) will req. transpose convolutions (for $\frac{\partial y_L}{\partial y_L}$) and regular convolutions (for $A_{w_L} y_{L-1}$ in σ' function).

5. (workings)

(a) Net 1



Net 2

