MFDNN

Homework 5

24/04/04

1. (Working)    (denoting loss function with $\mathcal{L}$)

$'dy' := \dfrac{\partial \mathcal{L}}{\partial y_L} = \dfrac{\partial}{\partial y_L}\left(\tfrac{1}{2}(f_\theta(x) - y)^2\right) = \dfrac{\partial}{\partial y_L}\left(\tfrac{1}{2}(y_L - y)^2\right) = y_L - y$   (as $g$ iven)

$z_\ell := A_\ell\, y_{\ell-1} + b_\ell$

For $\ell = L$ ('ell'$= L-1$),    $\dfrac{\partial \mathcal{L}}{\partial b_L} = \dfrac{\partial \mathcal{L}}{\partial y_L}\dfrac{\partial y_L}{\partial b_L} = dy \cdot 1 = dy \cdot \overset{\sigma'}{S(z)}$

For $\ell < L$ ('ell'$< L-1$),

$\dfrac{\partial \mathcal{L}}{\partial b_\ell} = \overbrace{\dfrac{\partial \mathcal{L}}{\partial y_L}\dfrac{\partial y_L}{\partial y_\ell}}^{=:'dy'}\dfrac{\partial y_\ell}{\partial b_\ell} = \dfrac{\partial \mathcal{L}}{\partial y_L}\overbrace{\left(\dfrac{\partial y_L}{\partial y_{L-1}}\cdots\dfrac{\partial y_{\ell+1}}{\partial y_\ell}\right)}\dfrac{\partial y_\ell}{\partial b_\ell} = dy \cdot \text{diag}(S(z))$

so 'ell' to $dy$ term to be reused next time

so $dy' := dy \cdot \text{diag}(\sigma'(z))A_\ell$

i.e. $dy := \dfrac{\partial \mathcal{L}}{\partial y_{\ell-1}} = \dfrac{\partial \mathcal{L}}{\partial y_L}\dfrac{\partial y_L}{\partial y_{\ell-1}}$

$\dfrac{\partial \mathcal{L}}{\partial A_\ell} = \dfrac{\partial \mathcal{L}}{\partial y_L}\dfrac{\partial y_L}{\partial A_\ell}$

✶   $= \dfrac{\partial \mathcal{L}}{\partial y_L}\,\text{diag}(S(z))\left(\dfrac{\partial y_L}{\partial y_\ell}\right)^T (y_{\ell-1})^T$

$$\underbrace{A_1, A_2, \ldots, A_\ell}, \overbrace{A_{\ell+1}, \ldots, A_L}^{A_j \text{ small}}$$

2. $\dfrac{\partial y_L}{\partial b_i} = \dfrac{\partial y_L}{\partial y_{L-1}} \dfrac{\partial y_{L-1}}{\partial y_{L-2}} \cdots \dfrac{\partial y_{i+1}}{\partial y_i} \dfrac{\partial y_i}{\partial b_i}$     by chain rule

$\qquad = f(A_i, \ldots, A_L) \, \text{diag}\left(\underline{\sigma'(A_i y_{i-1} + b_i)}\right)$     by HW 4.6

$\qquad\qquad\qquad\qquad\qquad\qquad$ where $f$ is a function involving the matrix multiplication of its args

$A_j \in \{A_i, \ldots, A_L\}$, i.e. $A_j$ is a term in $f$    $\longleftarrow$ inc. $\text{diag}(\sigma'(A_\ell y_{\ell-1} + b_\ell))$

If $A_j$ is small, then $f$ will be a matrix multiplication of not too large matrices and a small matrix, $A_j$, so the result will be small and hence $\dfrac{\partial y_L}{\partial b_i}$ become small.

$$\frac{\partial y_L}{\partial A_i} = \text{diag}\left(\underline{\sigma'(A_i y_{i-1} + b_i)}\right) \underbrace{\left(\frac{\partial y_L}{\partial y_{L-1}} \frac{\partial y_{L-1}}{\partial y_{L-2}} \cdots \frac{\partial y_{i+1}}{\partial y_i}\right)^T}_{f(A_i, \ldots, A_L)} (y_{i-1})^T$$

Here, by a similar argument as above, $\dfrac{\partial y_L}{\partial A_i}$ become small ($f$ will be small)

If $|\tilde{y}_j|$ is large then $\sigma'(\tilde{y}_j)$ is small (tending to 0) so now the $\text{diag}(\sigma'(\tilde{y}_i))$ term will be the small one in a matrix multiplication of otherwise not too large matrices and so $\dfrac{\partial y_L}{\partial b_i}, \dfrac{\partial y_L}{\partial A_i}$ become small.

3. For $k=0$,   $\theta_I^1 = \theta^0 - \alpha g^0 + \beta(\theta^0 - \theta^{-1})$

$\qquad\qquad = \theta^0 - \alpha g^0$     by Form I

$\qquad\quad \theta_{II}^1 = \theta^0 - \alpha v^1$

$\qquad\qquad = \theta^0 - \alpha(g^0 + \beta v^0)$

$\qquad\qquad = \theta^0 - \alpha g^0$     by Form II     So forms equivalent for $k=0$.

For $k=1$,   $\theta_I^2 = \theta^1 - \alpha g^1 + \beta(\theta^1 - \theta^0) = \theta^1 - \alpha g^1 + \beta(-\alpha v^1)$    since $\theta^1$ and $\theta^0$ equivalent

$\qquad\qquad\qquad\qquad = \theta^1 - \alpha g^1 - \alpha \beta v^1$

$\qquad\quad \theta_{II}^2 = \theta^1 - \alpha v^2$

$\qquad\qquad = \theta^1 - \alpha(g^1 + \beta v^1) = \theta^1 - \alpha g^1 - \alpha \beta v^1$    So equivalent for $k=1$.

Assume true for $k=n, n-1$ (i.e. $\theta_I^n = \theta_{II}^n = \theta^n$ and $\theta_I^{n-1} = \theta_{II}^{n-1} = \theta^{n-1}$)

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\longleftarrow$ no need to distinguish since assume equal

Consider $\theta^{n+1}$ term. By form I,   $\theta_I^{n+1} = \theta^n - \alpha g^n + \beta(\theta^n - \theta^{n-1})$

$\qquad\qquad\qquad\quad$ By form II,   $\theta_{II}^{n+1} = \theta^n - \alpha(g^n + \beta v^n)$

By assumption,   $\theta_I^{n+1} = \left(\theta_{II}^{n+1} + \alpha(g^n + \beta v^n)\right) - \alpha g^n + \beta(\theta^n - \theta^{n-1})$

$\qquad\qquad\qquad = \theta_{II}^{n+1} + \alpha \beta v^n + \beta(\theta^n - \theta^{n-1})$

$\qquad\qquad\qquad = \theta_{II}^{n+1} + \alpha \beta \left(\dfrac{\theta^n - \theta^{n-1}}{-\alpha}\right) + \beta(\theta^n - \theta^{n-1})$     since $\theta^k = \theta^{k-1} - \alpha v^k$ for $k \geq 1$
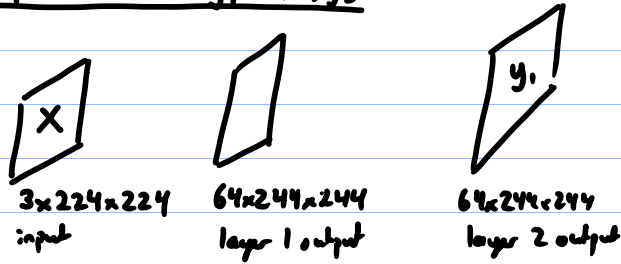
$\qquad\qquad\qquad = \theta_{II}^{n+1}$    so forms are equivalent for $k=n+1$. Since true for $k=0,1$, true for all $k$ by induction

4. Receptive field of $y_1[k,i,j]$

NB: using 0-indexing for $i,j,m,n$



3×224×224
input

64×244×244
layer 1 output
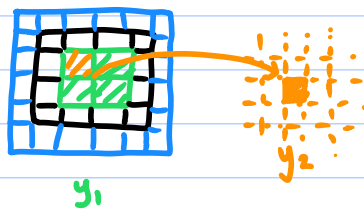
$y_1$
64×244×244
layer 2 output

After first layer is applied, each new 'pixel' depends on the 8 surrounding pixels of input in all 3 channels.
Applying another convolution each new 'pixel' depends again on the 8 surrounding it in all 64 channels.



$y_1[k,i,j]$
layer 1 output
input

So $y_1[k,i,j]$ depends on $X[c,m,n]$
for $1 \le c \le 3$, $i-2 \le m \le i+2$, $j-2 \le n \le j+2$
(just treat 'negative bounds' as 0 bounds)

## Receptive field of $y_2$

Applying a max pool layer changes the diagram slightly since each $y_2$ pixel depends on 4 $y_1$ pixels



$y_1$

So $y_2[k,i,j]$ depends on $X[c,m,n]$
for $1 \le c \le 3$, $2i-2 \le m \le 2i+3$, $2j-2 \le n \le 2j+3$

## Receptive field of $y_3$ – using recurrence relation equations

| | $k_i$ | $s_i$ | $p_i$ | |
|---|---|---|---|---|
| input, $y_0$ → $i=1$ | 3 | 1 | 1 | |
| $i=2$ | 3 | 1 | 1 | |
| $y_1$ → max $i=3$ | 2 | 2 | 0 | (L=6) |
| $y_2$ → $i=4$ | 3 | 1 | 1 | |
| $i=5$ | 3 | 1 | 1 | |
| max $i=6$ | 2 | 2 | 0 | |
| $y_3$ → | | | | |

Layer definitions ↑

$$r_0 = \sum_{l=1}^{L}\left((k_l-1)\prod_{i=1}^{l-1}s_i\right)+1 \quad \text{(equation 2 from reference)}$$

$= (2)(1) + (2)(1) + (1)(1\times1) + (2)(1\times1\times2) + (2)(1\times1\times2\times1) + (1)(1\times1\times2\times1\times1) + 1$

$= 2+2+1+4+4+2+1 = 16 \Rightarrow$ receptive field size of $y_3$ is $16\times16$
(also agrees with results for $y_1$ and $y_2$)

RF left-side index, $u_0 = u_L \prod_{n=1}^{L} s_n - \sum_{l=1}^{L} p_l \prod_{n=1}^{l-1} s_n \quad$ (equation 5)

(L=6)          or $j$, symmetric

$= (i)(1\times1\times2\times1\times1\times2) - \left((1+1(1)+0+1(1\times1\times2)+1(1\times1\times2\times1)+0\right)$

$= 4i - 6$

RF right-side index, $v_0 = v_L \prod_{n=1}^{L} s_n - \sum_{l=1}^{L}(1+p_l-k_l)\prod_{n=1}^{l-1}s_n \quad$ (equation 6)
(L=6)

$= (i)(1\times1\times2\times1\times1\times2) - \left(-1 + (-1)(1) + (-1)(1\times1) + (-1)(1\times1\times2) + (-1)(1\times1\times2\times1) + (-1)(1\times1\times2\times1\times1)\right)$

$= 4i - (-9) = 4i + 9$

So $y_3[k,i,j]$ depends on
$X[c,m,n]$ for $1 \le c \le 3$, $4i-6 \le m \le 4i+9$, $4j-6 \le n \le 4j+9$

**5.**

| | **Naïve Inception** | **With Bottlenecks** |

(i) Num trainable parameters

slide 12 $\left[k^2 C_{in} C_{out} + C_{out}\right]$

**Naïve Inception:**
$$(1^2 \times 256 \times 128 + 128) +$$
$$(3^2 \times 256 \times 192 + 192) +$$
$$(3^2 \times 256 \times 96 + 96)$$
$$= 696,736$$

**With Bottlenecks:**
$$(1^2 \times 256 \times 128 + 128) +$$
$$(1^2 \times 256 \times 64 + 64) +$$
$$(3^2 \times 64 \times 192 + 192) +$$
$$(1^2 \times 256 \times 64 + 64) +$$
$$(5^2 \times 64 \times 96 + 96) +$$
$$(1^2 \times 256 \times 64 + 64)$$
$$= 346,720$$

(ii) Each output element requires summing $k^2$ elements ($k$ is filter size), so $k^2 - 1$ additions, over $C_{in}$ layers so $C_{in}(k^2-1) + (C_{in}-1)$ additions. This is performed for each output layer and a bias is added, so: $C_{out}(C_{in}(k^2-1) + (C_{in}-1) + 1) = C_{out} C_{in} k^2 = C_{in} \times C_{out} \times k^2$. This is for each window location. So repeat for each 'pixel' of output i.e. $m \times n$ (since all padded to maintain dimensions)

slide 7 $\left[Y_{\ell, i, j} = \sum_{\gamma}^{C_{in}} \sum_{\alpha}^{k} \sum_{\beta}^{k} (\ldots) + b_\ell\right]$

$\ell = 1, \ldots, C_{out}$

**Additions:**

**Naïve Inception:**
$$((256 \times 128 \times 1^2) +$$
$$(256 \times 192 \times 3^2) +$$
$$(256 \times 96 \times 5^2)) \times 32^2$$
$$= 1,115,684,864$$

**With Bottlenecks:**
$$((256 \times 128 \times 1^2) +$$
$$(256 \times 64 \times 1^2) +$$
$$(64 \times 192 \times 3^2) +$$
$$(256 \times 64 \times 1^2) +$$
$$(64 \times 96 \times 5^2) +$$
$$(256 \times 64 \times 1^2)) \times 32^2$$
$$= 354,418,688$$

Multiplications are simpler to count since it's merely counting $(\ldots)$ term $\Rightarrow C_{out} C_{in} k^2$ (coincidentally equal!) So same num of multiplications as additions.

The activation function is evaluated on each element of a layer, so on $C_{out} \times m \times n$ elements.

**Naïve Inception:**
$$(128 \times 32^2) +$$
$$(192 \times 32^2) +$$
$$(96 \times 32^2)$$
$$= 425,984$$

**With Bottlenecks:**
$$(128 \times 32^2) +$$
$$(64 \times 32^2) +$$
$$(192 \times 32^2) +$$
$$(64 \times 32^2) +$$
$$(96 \times 32^2) +$$
$$(64 \times 32^2)$$
$$= 622,592$$

$$\sum_x^3 x = 1 + 2 + 3$$
(↑ ↑ 2 addition)

$$\sum_k^3 \sum_x^3 \sum_y^3 xy = (1+2+3) + (2+4+6) + (3+4+9) + \ldots + \ldots$$
(↑ 1/8 additions (4 per group)  ↑ 8 additions ↑ 8 additions)

k sets / terms within / 'joining' addition between sets
$$k(k-1) + (k-1)$$
$$= (k-1)(k+1) = (k^2 - 1)$$