# NN and Deep Learning

## Math Exercise 3

# Math 1

a) Since the convolution kernel can be characterized as

$$(x * W)_k = \sum_{i=1}^{W} w_i \, x_k \qquad \text{for } x \in \mathbb{R}^D.$$

Since there is no Padding we see that there are $D - W + 1$ valid placements for the kernel.

And since there are $M$ convolutions, there are total of $\underline{M * (D - W + 1)}$ outputs.

Because convolutions have weights, and there are $M$ convolutions, there are

$$\underline{M * W \quad \text{weights}}$$

6) In a fully connected layer there needs to be

$$M \times (W - D + 1) \text{ neurons}$$

to replicate the convolution behaviour.

Neurons would connect to all D inputs.

Thus, there are

$$M \times (W - D + 1) \times D \text{ weights},$$

most of which should be set to zero.

c) The fully connected layer can be characterized by

$$Z \in \mathbb{R}^{M(W-D+1) \times D}$$

matrix that maps $x \mapsto Zx =: z$.

In each neuron we want to have

$$z_k = \sum_{i=1}^{D} Z_{k,i} \, x_i$$

$$\approx \sum_{i=1}^{W} Z_{k, k-\lfloor \frac{W}{2} \rfloor + i} \, x_i \quad ,$$

meaning that entries outside the convolution window should be penalized and shrinked to zero.

We can achieve this by defining shrinkage matrix $P$ s.t.

$$P_{i,k} = \begin{cases} 1, & i \notin \{k - \lfloor \frac{W}{2} \rfloor + i \mid i \in \mathbb{N}, \, i \leq W\} \\ 0, & i \in =\!/\!/\!= . \end{cases}$$

Then regularization

$$R(x) := \| Px \|_1$$

penalizes those entries that are not within the convolution kernels window and the $\ell^1$ norm enforces sparsity.

B)

# Math 2

a) The blocks can be drawn as

1: $\quad x \xmapsto{\;f\;} f_1(x) \xmapsto{ReLU} ReLU(f_1(x)) \xmapsto{\oplus} x + ReLU(f_1(x))$

2: $\quad x \xmapsto{\;f_1\;} f_1(x) \xmapsto{ReLU} ReLU(f_1(x)) \xmapsto{\;f_2\;} f_2(ReLU(f_1(x))) \to \oplus\; f_1(x) + f_2 \dots$

$\qquad\qquad\qquad\qquad\qquad\qquad f_1 \qquad\qquad\qquad\qquad\qquad\qquad\qquad \to f_1(x)$

3

$\quad x \xmapsto{ReLU} ReLU(x) \xmapsto{\;f_1\;} f_1(ReLU(x)) \xmapsto{ReLU} ReLU(f_1(ReLU(x))) \xmapsto{\;f_2\;} f_2(ReLU(f_1(ReLU(x))))$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \to \oplus \mapsto y$

b)

1: 
$$y = \begin{cases} x & , f_1(x) = w_1 x + b_1 \leq 0 \\ (1+w_1)x + b_1 & , f_1(x) > 0 \end{cases}$$

Thus, $y$ is piecewise linear.

2:
$$y = \begin{cases} w_1 x + b_1 + b_2 & , f_1(x) \leq 0 \\ f_1(x) + w_2 f_1(x) + b_2 & \\ \quad = (1+w_2)(w_1 x + b_1) + b_2 & , f_1(x) > 0 \end{cases}$$

piecewise linear with knot at $f_1(x) = 0$.

3:
$$y = \begin{cases} x + f_2(ReLU(b_1)) = \begin{cases} w_1 x + b_2 & , x \leq 0, b_1 \leq 0 \\ f_1(x) + w_2 b_1 + b_2 & , x \leq 0, b_1 > 0 \end{cases} \\ x + f_2(ReLU(f_1(x))) = \begin{cases} x + f_2(b) & , f_1(x) \leq 0, x > 0 \\ x + f_2(f_1(x)) & , f_1(x) > 0, x > 0 \end{cases} \end{cases}$$

c) Note the derivative is not actual at knots. Where the function is differentiable we get:

1:

$$\frac{dy}{dx} = \begin{cases} 1 & , \; f_1(x) \leq 0 \\ 1 + w_2 & , \; f_1(x) > 0 \end{cases}$$

2:

$$\frac{dy}{dx} = \begin{cases} w_1 & , \; f_1(x) \leq 0 \\ (1 + w_2) w_0 & , \; f_1(x) \geq 0 \end{cases}$$

3:

$$\frac{dy}{dx} = \begin{cases} w_1 & , \; x \leq 0 \\ 1 & , \; x > 0, \; f_1(x) \leq 0 \\ 1 + w_1 w_2 x & , \; x > 0, \; f_1(x) > 0. \end{cases}$$

would probably prefer block 1, because it can be used to produce all the same functions as the more complicated blocks.

Block 1 also has the "most" stable gradient with block two.