

Homework 5Exercise 1:

The NAND gate behaves as follows

Input		Output
$x_1$	$x_2$	$y$
0	0	1
1	0	1
0	1	1
1	1	0

let  $(x_1, x_2)$  be the input to a perceptron

$w_1, w_2, b$  the weights and bias of the perceptron.

Then  $y = f(x_1, x_2) = \begin{cases} 1 & \text{if } w_1 x_1 + w_2 x_2 + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$

is the output of the perceptron.

In that case, we want:

- $f(0,0) = f(1,0) = f(0,1) = 1$
- $f(1,1) = 0$

We can choose

$$w_1 = w_2 = -2 \quad b = 3$$

Exercise 1.2:

let  $x$  the input image (grayscale)

$K$  the kernel of a convolution layer, of size  $n$ .

$y$  the output of  $x * K$ .

We apply the convolution in 1D.

$$y[k] = \sum_{j=-n} x[j] K[k-j]$$

Let's write a DCT<sub>I</sub> transform in 1D as

$$Y_k = \frac{1}{\sqrt{N}} \sum_{j=0}^{N-1} x_j \cos \left[ \pi \left( j + \frac{1}{2} \right) \frac{k}{N} \right]$$

with  $N=4$  on a  $4 \times 4$  patch, in 1D

and  $k \leq N$



Thus, we can create  $N$  kernels:  $K_k$ ,  $k \leq N$  such that

$$Y_k = 2\alpha_k \sum_{j=0}^{N-1} X_j \cos \left[ \pi \left( j + \frac{1}{2} \right) \frac{k}{N} \right]$$

$$= \sum_{j=0}^{N-1} X_j \underbrace{2\alpha_k \cos \left[ \pi \left( -(k-j) + k + \frac{1}{2} \right) \frac{k}{N} \right]}_{\text{red arrow } \leftarrow k-k}$$

$$\text{ie } Y[k] = \sum_{j=0}^{N-1} X[j] K_k[k-j]$$

Thus, we can get  $N$  filters with the  $k$ -th filter:

$$K_k[i] = 2\alpha_k \cos \left[ \pi \left( -i + k + \frac{1}{2} \right) \frac{k}{N} \right]$$

and  $0 \leq i \leq N-1$ .

In the end, we have 4 filters of size 4 each.

The DCT transform has been represented as a convolution in 1D.

If we do the same process in 2D, we represent DCT as 2 convolutions.

### Exercice 13

let  $x \in \mathbb{R}^3$

$W \in \mathbb{R}^{4 \times 3}$

$b \in \mathbb{R}^4 = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$

$y = f(x) = g(Wx + b) \in \mathbb{R}^4$

$g$  the sigmoid function defined as  $g: u \mapsto \frac{1}{1+e^{-u}}$

Then, let's compute the gradient of  $f$  wrt  $W, b, x$ .

Preliminary result:

let's call  $\tilde{y} = Wx + b$ . Then, for  $k_j = \begin{bmatrix} i \\ j \end{bmatrix} \in \mathbb{R}^2$  (red arrow  $\leftarrow i, j$ -th)

•  $(w_i + b_j)^T x - w_i^T x = b_j^T x = x^T b_j$

So  $\frac{\partial \tilde{y}}{\partial W_{ij}} = x^T$

• Similarly,  $\frac{\partial \tilde{y}}{\partial b^*} = 1$



$$\bullet \frac{\partial \tilde{y}}{\partial x_i} = W^T$$

And

$$\frac{\partial y}{\partial \tilde{y}} = \frac{e^{-\tilde{y}}}{(1 + e^{-\tilde{y}})^2} = \frac{1}{e^{\tilde{y}} + 2 + e^{-\tilde{y}}}$$

And finally, using the chain rule:

$$\bullet \frac{\partial f}{\partial W} = \frac{\partial f}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial W} = \frac{x}{2 + e^{-\tilde{y}} + e^{\tilde{y}}}$$

$$\bullet \frac{\partial f}{\partial x} = \frac{\partial f}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial x} = \frac{W^T}{2 + e^{-\tilde{y}} + e^{\tilde{y}}}$$

$$\bullet \frac{\partial f}{\partial b} = \frac{\partial f}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial b} = \frac{1}{2 + e^{-\tilde{y}} + e^{\tilde{y}}}$$

To be more precise, we have

$$\bullet \frac{\partial f_k}{\partial W_{ij}} = \frac{\partial f_k}{\partial \tilde{y}_k} \times \frac{\partial \tilde{y}_k}{\partial W_{ij}} = \begin{cases} 0 & \text{if } k \neq i \\ \frac{x_j}{2 + e^{-\tilde{y}_k} + e^{\tilde{y}_k}} & \text{otherwise} \end{cases}$$

$$\bullet \frac{\partial f_k}{\partial x_i} = \frac{\partial f_k}{\partial \tilde{y}_k} \times \frac{\partial \tilde{y}_k}{\partial x_i}$$

$$\bullet \frac{\partial f_k}{\partial x_i} = \frac{W_{ki}}{2 + e^{-\tilde{y}_k} + e^{\tilde{y}_k}}$$

$$\bullet \frac{\partial f_k}{\partial b_i} = \frac{1}{2 + e^{-\tilde{y}_k} + e^{\tilde{y}_k}}$$



#### Exercise 1.4:

Let  $f_i(x, \theta_i)$  a network layer

$F(x)$  the network of 3  $f_i$  layers

$G(x)$  the network of 3  $f_i$  layers + a skip connection in last layer.

Let's compute the gradient, using the chain rule:

We set  $\tilde{y} = f_1(x, \theta_1)$

Thus

$$\bullet \frac{\partial F}{\partial \theta_1} = \frac{\partial f_3(y, \theta_3)}{\partial y} \times \frac{\partial f_2(\tilde{y}, \theta_2)}{\partial \tilde{y}} \times \frac{\partial f_1(x, \theta_1)}{\partial \theta_1}$$

$$\bullet \frac{\partial F}{\partial \theta_2} = \frac{\partial f_3(y, \theta_3)}{\partial y} \times \frac{\partial f_2(\tilde{y}, \theta_2)}{\partial \theta_2}$$

$$\bullet \frac{\partial F}{\partial \theta_3} = \frac{\partial f_3(y, \theta_3)}{\partial \theta_3}$$

If we do the same on  $G$ :

$$\bullet \frac{\partial G}{\partial \theta_1} = \left[ 1 + \frac{\partial f_3(y, \theta_3)}{\partial y} \right] \times \frac{\partial f_2(\tilde{y}, \theta_2)}{\partial \tilde{y}} \times \frac{\partial f_1(x, \theta_1)}{\partial \theta_1}$$

$$\text{ie } \frac{\partial G}{\partial \theta_1} = \frac{\partial F}{\partial \theta_1} + \frac{\partial f_2(\tilde{y}, \theta_2)}{\partial \tilde{y}} \times \frac{\partial f_1(x, \theta_1)}{\partial \theta_1}$$

$$\bullet \frac{\partial G}{\partial \theta_2} = \frac{\partial F}{\partial \theta_2} + \frac{\partial f_2(\tilde{y}, \theta_2)}{\partial \theta_2}$$

$$\bullet \frac{\partial G}{\partial \theta_3} = \frac{\partial F}{\partial \theta_3}$$

Adding a skip connection prevents vanishing gradients as the gradient becomes a sum of multiple terms.