

Exercise 1.1

Implement NAND gate with a single perceptron.
We have definition of NAND:

A	B	ANANDB
0	0	1
0	1	1
1	0	1
1	1	0

And the perceptron:
 $f(x) = \begin{cases} 1, & \text{if } wx + b > 0 \\ 0, & \text{otherwise} \end{cases}$

In our case:
 $f((A, B)) = \begin{cases} 1, & \text{if } w_1 A + w_2 B + b > 0 \\ 0, & \text{otherwise} \end{cases}$

We can choose $w = (-0.5, -0.5)$ and $b = 1$. Then, when $w \neq (1, 1)$ $wx + b > 0$ and when $w = (1, 1)$ $wx + b = 0$.

Exercise 1.2 Patch-wise DCT transform as convolution.
Describe how the patch-wise DCT transform (with 4×4 patches) of a grayscale image can be implemented and represented using convolutions.

We have DCT for grayscale image:
 $y_k = 2d_k \sum_{j=0}^{N-1} x_j \cos(\pi(j+\frac{1}{2})\frac{k}{N})$, $d_k = \begin{cases} \frac{1}{\sqrt{4N}}, & k=0 \\ \frac{1}{\sqrt{2N}}, & \text{otherwise} \end{cases}$
In our case, we have patch $4 \times 4 \Rightarrow N=4$.

We have convolution definition $N=7$
 $y[k] = (x * z)[k] = \sum_{j=0}^{N-1} x[j] z[k-j]$

Let's perform transformation:

$$y_k = 2d_k \sum_{j=0}^3 x_j \cos(\pi(j+\frac{1}{2})\frac{k}{4}) = \sum_{j=0}^3 x_j \cdot (2d_k \cos(\pi(j+\frac{1}{2})\frac{k}{4}))$$

We can say that $x_j = xz_j$ and second part is $z_k = z_j$.
 Therefore, $y_k = xz_k$, $z_k = (z_{kj})$, $z_{kj} = z_k \cos(\pi(j-\frac{1}{2})\frac{k}{N})$,
 $x = (x_j)_{0 \leq j \leq 3}$, $z_k = \begin{cases} \frac{1}{\sqrt{2}}, & k=0 \\ \frac{1}{\sqrt{2}}, & \text{otherwise} \end{cases}$

Exercise 1.3

Use the chain rule to compute the gradient of $f(x) = g(Wx+b)$ with respect to W, x, b , where g is sigmoid and W is 4×3 matrix.

We have next functions: $f(x) = (g \circ h)(x)$, $h(x) = Wx+b$,

$$g(y) = \frac{1}{1+e^{-y}}, \quad \frac{\partial f}{\partial W} = \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial W}$$

$$\frac{\partial f}{\partial h} = \frac{e^{-h}}{(1+e^{-h})^2}, \quad \frac{\partial h}{\partial W} = x^T \Rightarrow \frac{\partial f}{\partial W} = \frac{e^{-(Wx+b)}}{(1+e^{-(Wx+b)})^2} x^T$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial b}, \quad \frac{\partial h}{\partial b} = 1 \Rightarrow \frac{\partial f}{\partial b} = \frac{e^{-(Wx+b)}}{(1+e^{-(Wx+b)})^2}$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial x}, \quad \frac{\partial h}{\partial x} = W^T \Rightarrow \frac{\partial f}{\partial x} = \frac{e^{-(Wx+b)}}{(1+e^{-(Wx+b)})^2} W^T$$

$$\frac{\partial f}{\partial h} = \frac{e^{-h}}{(1+e^{-h})^2}, \quad \frac{\partial h}{\partial x} = W^T \Rightarrow \frac{\partial f}{\partial x} = \frac{e^{-(Wx+b)}}{(1+e^{-(Wx+b)})^2} W^T$$

Exercise 1.4

Let $f_i(x; \theta_i)$ denote a network layer, with parameter vector θ_i and inputs x . Consider the multi-layer network defined as $F(x) = f_3(y, \theta_3)$, where $y = f_2(f_1(x, \theta_1), \theta_2)$. Also, consider the multi-layer network with skip connection defined as $G(x) = y + f_3(y, \theta_3)$. Use the chain rule to compute the gradient of $F(x)$ with

respect to all θ_i . Do the same for $G(2)$. What is special in

$$F(x): \frac{\partial F}{\partial \theta_1} = \frac{\partial f_3}{\partial \theta_1}, \quad \frac{\partial F}{\partial \theta_2} = \frac{\partial f_3}{\partial \theta_2} \cdot \frac{\partial f_2}{\partial \theta_2}, \quad \frac{\partial F}{\partial \theta_1} = \frac{\partial f_3}{\partial \theta_2} \cdot \frac{\partial f_2}{\partial \theta_1} \cdot \frac{\partial f_1}{\partial \theta_1}$$

$$G(x): \frac{\partial G}{\partial \theta_3} = \frac{\partial (y + f_3(y, \theta_3))}{\partial \theta_3} = \frac{\partial y}{\partial \theta_3} + \frac{\partial f_3}{\partial \theta_3} = \frac{\partial f_3}{\partial \theta_3}, \text{ because } y \text{ is inde.}$$

$$\text{pendable from } \theta_3: \frac{\partial G}{\partial \theta_2} = \frac{\partial y}{\partial \theta_2} + \frac{\partial f_3}{\partial \theta_2} = \frac{\partial f_2}{\partial \theta_2} + \frac{\partial f_3}{\partial \theta_2} = \frac{\partial f_2}{\partial \theta_2} + \frac{\partial f_3}{\partial \theta_2} \cdot \frac{\partial f_2}{\partial \theta_2} =$$

$$= \frac{\partial f_2}{\partial \theta_2} (1 + \frac{\partial f_3}{\partial \theta_2})$$

$$\frac{\partial G}{\partial \theta_1} = \frac{\partial y}{\partial \theta_1} + \frac{\partial f_3}{\partial \theta_1} = \frac{\partial f_2}{\partial \theta_1} + \frac{\partial f_3}{\partial \theta_1} = \frac{\partial f_2}{\partial \theta_1} \cdot \frac{\partial f_1}{\partial \theta_1} + \frac{\partial f_3}{\partial \theta_2} \cdot \frac{\partial f_2}{\partial \theta_1} \cdot \frac{\partial f_1}{\partial \theta_1}$$

$$= (\frac{\partial f_2}{\partial \theta_1} \cdot \frac{\partial f_1}{\partial \theta_1}) (1 + \frac{\partial f_3}{\partial \theta_2})$$

The closer ~~output~~ calculations to the end, the ~~more~~ ^{more} severe the vanishing gradient problem. In $\frac{\partial G}{\partial \theta_1}$, we have $\frac{\partial f_3}{\partial \theta_2}$ which at some point can approach 0. However, $\frac{\partial G}{\partial \theta_1}$ won't approach 0 in that case because we have $(1 + \frac{\partial f_3}{\partial \theta_2})$. Therefore, $\frac{\partial G}{\partial \theta_1} \rightarrow 0$ as $\frac{\partial f_2}{\partial \theta_2} \rightarrow 0$, but $\frac{\partial G}{\partial \theta_1} \rightarrow \frac{\partial f_2}{\partial \theta_1} \cdot \frac{\partial f_1}{\partial \theta_1}$ which is ~~not~~ ^{is} the feature of skip connections: attenuate vanishing gradient problem.