# NNDL - 2025

## Exercise 1

## Prob. 1

C) Let $x \in \mathbb{R}$. Then, it $x \geq 0$, $\psi(x) = x = \max(\alpha x, x)$, because $\alpha x \leq x$ for $\alpha \in [0,1]$. Similarly, it $x < 0$, then $\alpha x \geq x$ and $\psi(x) = \alpha x = \max(\alpha x, x)$. Thus,

$$\psi(x) = \begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases} = \max(\alpha x, x).$$

a) Choose $\alpha = 0$, then

$$\psi(x) = \max(x, \alpha x) = \max(x, 0),$$

which shows that ReLu can be approximated by leaky ReLu, making it a special case.

b) Choose $\alpha := 1$. Then

$$\psi(x) = \max\{x, \alpha x\}$$

$$= \max\{x, x\} = x,$$

for every $x \in \mathbb{R}$

## Prob. 2

a) If $W_1 \in \mathbb{R}^{m \times n}$, i.e. first layer has $n$ neurons, then $W_2 \in \mathbb{R}^{n \times k}$ for some $k \in \mathbb{N}$. Thus, $W_2$ has $n$ rows.

b) If $NN$ is invertible (i.e. the product is injective), then $m \geq n$ meaning that $\dim R(A_{n+1}) \geq \dim R(A_n)$.

Assume $m < n$. Then there exists $x \neq 0$ s.t. $W_1 x = 0$, which implies

$$W_k \ldots W_1 x = 0 \quad \text{by linearity.}$$

Which would be a contradiction with injectivity since $W_k \ldots W_1 0 = 0$.

⨳

c) If all matrices are square then, the $NN$ is invertible iff all matrices are invertible.

Optimization:

## Prob 1:

Partial derivatives are given as

$$\partial_i f_1(w) = \partial_i \|w\|^2 = \partial_i \sum w_k^2$$

$$= \sum \partial_i w_u^2 = 2w_i.$$

Thus, the gradient is

$$\nabla f_1(w) = \begin{bmatrix} \partial_1 \\ \vdots \\ \partial_n \end{bmatrix} = 2 \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

## Prob 2

Hessian is given as

$$H_f = (D \nabla f)^T = \begin{bmatrix} D \nabla f_1 \\ \vdots \\ D \nabla f_2 \end{bmatrix}^T,$$

where $D$ is the differentiation operator.

Since $(D \nabla f_i)_j = \partial_j \partial_i f = \partial_j 2w_i = \begin{cases} 0, & j \neq i \\ 2, & j = i \end{cases}$

we set

$$H_f = 2I_{n\times n} \quad (\text{identity matrix}).$$

□

## Prob 3.

The Newton's method is given by the recursion

$$w_{k+1} = w_u - H_f(w_u)^{-1} \nabla f(w_u)$$

which in our case is

$$w_{u+1} = w_u - w_u$$

always converging in one step.

# Prob 4

The gradient is given as

$$\nabla f_0(w) = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} = z$$

# Prob 5

Using chain rule we obtain gradient at point $w$ as

$$\nabla f_z(w) = \left( Dg(w^T z) \, D(w^T z) \right)^T$$

$$= g'(w^T z) \, \nabla f_0(w)$$

$$= g'(w^T z) \, z$$

□

## Prob 6

Stochastic gradient is given as

$$\nabla f_3(w) = \nabla \mathbb{E}(s(w^T z))$$

$$= \mathbb{E}(\nabla s(w^T z)) \qquad \text{/By definition}$$

$$= \mathbb{E}(s'(w^T z) z)$$

$$\boxed{B}$$

## Prob 7

We note partial derivative of the quadratic form at point $w$ $\beta$ given as

$$\partial_i = \frac{1}{h}\left(((w + h e_i)^T A (w + h e_i) - w^T A w\right)$$

$$= \frac{1}{h}\left(w^T A h e_i + (h e_i)^T A w\right)$$

$$= \frac{1}{h}\left(w^T (A + A^T) e_i h\right)$$

$$= 2 w^T A e_i = 2 w^T A_{\cdot i}$$

And by definition the gradient is
$$\nabla \tfrac{1}{2} w^T A w = w^T A$$

$$\boxed{2}$$

## Prob 7

b) First more that $\|w\|^4 = (\|w\|^2)^2$. Then, by chain rule and solution 9., we get

$$D\|w\|^4 = D(\|w\|^2)^2 \, D\|w\|^2$$

$$= 2\|w\|^2 (2 \nabla \|w\|^2)^T$$

Thus so sing to get the gradient we have

$$D \frac{1}{4} \|w\|^4 = \|w\|^2 \, w$$