

Homework Week 4

MATHEMATICS OF DEEP LEARNING
MASH & IASD 2026

Lecturer: Bruno Loureiro, bruno.loureiro@di.ens.fr

Instructions: This homework is **due on Monday 23/02/2026**. Please upload your solutions in a PDF file named HW4_NOM_PRENOM.PDF [here](#). Formats accepted: PDF (LaTeX or a **readable** scan of handwritten solutions).

1 Exercises

Exercise 1.

Consider a two-layer neural network with ReLU activation $\sigma(x) = x_+$:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{\sqrt{p}} \sum_{j=1}^p a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle). \quad (1)$$

Assume that the weights are initialised as $a_j^0 \sim \text{Unif}(\{-1, 1\})$, $\mathbf{w}^0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

- (a) Show that the NTK kernel is given by:

$$\begin{aligned} K_{\text{NTK}}(\mathbf{x}, \mathbf{x}') &:= \langle \mathbf{x}, \mathbf{x}' \rangle \mathbb{E}_{a, \mathbf{w}} [a^2 \sigma'(\langle \mathbf{w}, \mathbf{x} \rangle) \sigma'(\langle \mathbf{w}, \mathbf{x}' \rangle)] \\ &= \langle \mathbf{x}, \mathbf{x}' \rangle \left[\frac{1}{2} - \frac{1}{2\pi} \arccos \left(\frac{\langle \mathbf{x}, \mathbf{x}' \rangle}{\|\mathbf{x}\|_2 \cdot \|\mathbf{x}'\|_2} \right) \right] \end{aligned} \quad (2)$$

- (b) Let $\mathbf{x}_i \in \mathbb{R}^d$ denote a batch of n independently sampled covariates, and assume $\mathbf{x}_i \in B(\mathbf{0}, 1)$. Using Hoeffding's inequality, show that if $p \geq \Omega(\epsilon^{-2} n^2 \log n / \delta)$, then with probability at least $1 - \delta$ over the random initialisation we have:

$$\|\hat{\mathbf{K}}_{\text{NTK}} - \mathbf{K}_{\text{NTK}}\|_{\text{F}} \leq \epsilon \quad (3)$$

where $\hat{\mathbf{K}}_{\text{NTK}}, \mathbf{K}_{\text{NTK}} \in \mathbb{R}^{n \times n}$ with:

$$\begin{aligned} \hat{\mathbf{K}}_{\text{NTK}, ij} &= \frac{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}{p} \sum_{k=1}^p a_k^2 \sigma'(\langle \mathbf{w}_k, \mathbf{x}_i \rangle) \sigma'(\langle \mathbf{w}_k, \mathbf{x}_j \rangle), \\ \mathbf{K}_{\text{NTK}, ij} &= K_{\text{NTK}}(\mathbf{x}_i, \mathbf{x}_j) \end{aligned} \quad (4)$$

- (c) Conclude that for large enough width p , we have that $\lambda_{\min}(\hat{\mathbf{K}}_{\text{NTK}}) > 0$ with high-probability.

Exercise 2.

Let $g(\mathbf{x}; \boldsymbol{\theta}) = \phi(f(\mathbf{x}; \boldsymbol{\theta}))$ where ϕ is a twice differentiable function and $f(\mathbf{x}; \boldsymbol{\theta})$ a two-layer neural network:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{\sqrt{p}} \sum_{j=1}^p a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle) \quad (5)$$

with twice-differentiable activation function σ .

- (a) Considering a_j to be fixed, show that for any $\boldsymbol{\theta}$, the Hessian matrix \mathbf{H}_g of g can be related to the Hessian matrix $\mathbf{H}(\boldsymbol{\theta})$ of f by:

$$\mathbf{H}_g(\boldsymbol{\theta}) = \phi'(f(\mathbf{x}; \boldsymbol{\theta}))\mathbf{H}(\boldsymbol{\theta}) + \phi''(f(\mathbf{x}; \boldsymbol{\theta}))\nabla_{\mathbf{w}}f(\mathbf{x}; \boldsymbol{\theta})\nabla_{\mathbf{w}}f(\mathbf{x}; \boldsymbol{\theta})^\top \quad (6)$$

- (b) Under standard initialisation $a^0 \sim \text{Unif}([-1, 1])$, what is the scaling in p of the operator norm of each of the terms above?
- (c) Conclude that $g(\mathbf{x}; \boldsymbol{\theta})$ does not linearise as $p \rightarrow \infty$.

Exercise 3.

Generalise the argument leading to Proposition 1 to the case where the second layer weights a_j are not fixed.