

Exam on Neural Network Theory February 11, 2021

Please note:

- Exam duration: 180 minutes
- Maximum number of points: 100
- You are not allowed to use any printed or handwritten material (i.e., books, lecture and discussion session notes, summaries), computers, tablets, smart phones or other electronic devices.
- Your solutions should be explained in detail and your handwriting needs to be clean and legible.
- Please do not use red or green pens. You may use pencils.
- Please note that the "ETH Zurich Disciplinary Code" (RS 361.1) applies.

Before you start:

- 1. The problem statements consist of 6 pages including this page. Please verify that you have received all 6 pages.
- 2. Please fill in your name, student ID card number and signature below.
- 3. Please place your student ID card at the front of your desk so we can verify your identity.

During the exam:

- 4. For your solutions, please use only the empty sheets provided by us. Should you need additional sheets, please let us know.
- 5. Each problem consists of several subproblems. If you do not provide a solution to a subproblem, you may, whenever applicable, nonetheless assume its conclusion in the ensuing subproblems.

After the exam:

- 6. Please write your name on every solution sheet and prepare all sheets in a pile. All sheets, including those containing problem statements, must be handed in.
- 7. Please clean up your desk and stay seated and silent until you are allowed to leave the room in a staggered manner row by row.
- 8. Please avoid crowding and leave the building by the most direct route.

Family name:	First name:
Student ID card No.:	
Signature:	

Problem 1 (25 points)

- (a) Realize the function $f(x) = 10|x|, x \in \mathbb{R}$, through a ReLU network Ψ (see the Handout for the definition of a ReLU network), with $W(\Psi) = 2$ and $\mathcal{B}(\Psi) = 2$. Specify the depth and the connectivity of the network Ψ .
- (b) Consider the ReLU network given by $\Phi : \mathbb{R} \to \mathbb{R}$

$$\Phi(x) = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \circ \rho \circ \left(\begin{pmatrix} 1 \\ 1 \\ 2.5 \end{pmatrix} x - \begin{pmatrix} 0 \\ 0.4 \\ 2 \end{pmatrix} \right).$$

Determine the depth, width, and connectivity of the network Φ according to the definition of ReLU networks specified in the Handout. Plot the function $\Phi(x)$ on the interval [-1, 1].

(c) Consider the following classification problem in \mathbb{R}^2 . The points $\{x_1, x_2, x_3\}$ in red constitute class 1 and the points $\{x_4, x_5, x_6\}$ in blue belong to class 2:



Show that the symmetric (around x = 0.5) function f(x), such that $f(x) = \Phi(2x)$ for $x \le 0.5$, with $\Phi(x)$ the ReLU network from subproblem (b), separates class 1 from class 2.¹ Realize f(x) through a depth-2 ReLU network. Specify the width and the connectivity of the resulting network.

(d) Recall the sawtooth function $g : [0,1] \rightarrow [0,1], g(x) = g_1(x) = \rho(2x) - \rho(4x - 2) + \rho(2x - 2)$ and its *s*-fold composition according to $g_s := g \circ g \circ \cdots \circ g, s \ge 2$.

Determine the cardinality of the set $\{x : g_s(x) = 1, x \in [0, 1]\}$ as a function of *s*.

(e) In subproblem (c) the class-separating function f(x) was realized by a ReLU network of depth 2. Using the sawtooth function from subproblem (d), find a deeper ReLU network of width 3 that realizes the class-separating function f(x). Specify the depth and the connectivity of this network.

¹We say that a function f(x) separates the given point set if all red points lie below the function graph and all blue points above.

Problem 2 (25 points)

- (a) Let (\mathcal{X}, ρ) be a metric space and \mathcal{C} a compact set in \mathcal{X} .
 - (i) State the definitions of ϵ -covering of C and of the ϵ -covering number $N(\epsilon; C, \rho)$.
 - (ii) State the definitions of ϵ -packing of C and of the ϵ -packing number $M(\epsilon; C, \rho)$.
 - (iii) Order the following quantities

 $N(\epsilon; \mathcal{C}, \rho), N(2\epsilon; \mathcal{C}, \rho), M(\epsilon; \mathcal{C}, \rho), M(2\epsilon; \mathcal{C}, \rho).$

(b) Consider the following parametric class of functions

$$\mathcal{F} = \{ f_{\theta} : [0, 1] \to \mathbb{R} \mid \theta \in [0, 1] \},\$$

where for any fixed $\theta \in [0, 1]$, we set $f_{\theta}(x) := \ln(1 + \theta x), x \in [0, 1]$. We take the ∞ -norm of functions defined on [0, 1] to be given by

$$||f - g||_{\infty} = \max_{x \in [0,1]} |f(x) - g(x)|.$$

(i) For $\epsilon < 1/2$, construct an ϵ -covering $A(\epsilon)$ for the class \mathcal{F} as follows. Set $T = \lfloor \frac{1}{2\epsilon} \rfloor$, and for $i = 0, 1, \ldots, T$, define $\theta_i = 2\epsilon i$. We also add the point $\theta_{T+1} = 1$, thereby forming a collection $\{\theta_0, \ldots, \theta_T, \theta_{T+1}\}$ contained in [0, 1]. Show that the associated functions $\{f_{\theta_0}, \ldots, f_{\theta_T}, f_{\theta_{T+1}}\}$ constitute an ϵ -covering of \mathcal{F} . Find an upper bound on the ϵ -covering number $N(\epsilon; \mathcal{F}, \|\cdot\|_{\infty})$ in terms of ϵ .

Hint: You can use without proof that $\frac{x}{1+x} \leq \ln(1+x) \leq x$, for all x > -1.

(ii) For $\epsilon < 1/3$, construct an ϵ -packing $P(\epsilon)$ for the class \mathcal{F} . Find a lower bound on the ϵ -packing number $M(\epsilon; \mathcal{F}, \|\cdot\|_{\infty})$ in terms of ϵ .

Hint: You can use without proof that $\frac{x}{1+x} \leq \ln(1+x) \leq x$, for all x > -1.

(iii) Show that the metric entropy of the class \mathcal{F} w.r.t. the norm $\|\cdot\|_{\infty}$ satisfies

$$\log N(\epsilon; \mathcal{F}, \|\cdot\|_{\infty}) \asymp \log(1/\epsilon), \text{ as } \epsilon \to 0.^2$$

²One writes $f \asymp g$, if f = O(g) and g = O(f). One writes f = O(g), if $\limsup_{\epsilon \to 0} \left| \frac{f(\epsilon)}{g(\epsilon)} \right| < \infty$.

Problem 3 (25 points)

Suppose $d, N \in \mathbb{N}$. We use the notation $||A||_{\infty}$ and $||y||_{\infty}$ for the maximum absolute value of the entries of matrix A and the components of vector y, respectively. Consider a single-hidden-layer logistic neural network with d-dimensional input and scalar output given by

$$\Phi(x) := A_2 \sigma(A_1 x), \quad x \in [-1, 1]^d,$$

where A_1 is an $N \times d$ matrix, A_2 is a $1 \times N$ matrix such that $||A_1||_{\infty}$, $||A_2||_{\infty} \leq 1$, and the logistic activation function given by $\sigma(x) = e^x/(1+e^x)$, $x \in \mathbb{R}$, acts componentwise, i.e., $\sigma(x_1, x_2, \ldots, x_N) = (\sigma(x_1), \sigma(x_2), \ldots, \sigma(x_N))$. In practice, the entries in A_1 and A_2 can not be stored with infinite precision in a computer. To take this constraint into account, we fix a natural number b and introduce a truncation operator that maps real numbers to their storable approximations as follows. For the real number a, we consider its binary representation truncated after the b-th digit after the comma and we define $T_b(a)$ to be the number corresponding to this truncated representation. For example, if a has binary representation $101, 10 \dots 10$ 1001, then $T_b(a)$ has binary representation $101, 10 \dots 10$. For matrices and vectors, the operator T_b acts on each entry, i.e., for an M-by-N matrix $A = (a_{i,j})_{1 \leq i \leq M}, M, N \in \mathbb{N}$, we have $T_b(A) = (T_b(a_{i,j}))_{1 \leq i \leq M}, M_b \in \mathbb{N}$.

for an *M*-by-*N* matrix $A = (a_{i,j})_{\substack{1 \le i \le M \\ 1 \le j \le N}}$, $M, N \in \mathbb{N}$, we have $T_b(A) = (T_b(a_{i,j}))_{\substack{1 \le i \le M \\ 1 \le j \le N}}$ likewise for vectors. Let $\tilde{A}_i = T_b(A_i)$, i = 1, 2, and let

$$\tilde{\Phi}(x) := \tilde{A}_2 \,\sigma\big(\tilde{A}_1 x\big), \quad x \in [-1, 1]^d,$$

be the network actually implemented. The goal of this problem is to derive an upper bound on the approximation error

$$\sup_{x \in [-1,1]^d} \left| \Phi(x) - \tilde{\Phi}(x) \right|$$

incurred by truncation for finite precision as just described.

(a) Show that for $a \in \mathbb{R}$, the error incurred by truncating to b bits after the comma is upper-bounded by 2^{-b} , i.e.,

$$|T_b(a) - a| \le 2^{-b}$$
, for all $a \in \mathbb{R}$.

(b) Show that

$$\sup_{x \in [-1,1]^d} \left\| A_1 x - \tilde{A}_1 x \right\|_{\infty} \le d \, 2^{-b}.$$

Hint: Use the inequality $||Ax||_{\infty} \leq n ||A||_{\infty} ||x||_{\infty}$, for all $m \times n$ matrices A and all n-vectors $x, m, n \in \mathbb{N}$.

(c) In order to understand the effect of the activation function σ in the propagation of the truncation error across layers in the network, we first need to establish

some basic properties of σ . Show that σ is Lipschitz-continuous with Lipschitz constant $\frac{1}{4}$ and that it takes value in [0, 1], i.e.,

$$\begin{aligned} |\sigma(x) - \sigma(y)| &\leq \frac{1}{4} |x - y|, \quad \text{for all } x, y \in \mathbb{R}, \\ 0 &\leq \sigma(x) \leq 1, \quad \text{for all } x \in \mathbb{R}. \end{aligned}$$

Hint: Compute the first and the second derivative of σ , find the range of values σ' takes on, and use the mean value theorem, which states that for every $x, y \in \mathbb{R}$ with $x \leq y$, there exists a real number z such that $x \leq z \leq y$ and

$$\sigma(x) - \sigma(y) = \sigma'(z)(x - y).$$

(d) Show that

$$\sup_{x\in[-1,1]^d} \left\| \sigma(A_1x) - \sigma(\tilde{A}_1x) \right\|_{\infty} \leq \frac{1}{4} d \, 2^{-b}.$$

Hint: Use Lipschitz continuity of σ as established in subproblem (c).

(e) Show that

$$\sup_{x \in [-1,1]^d} \left| \Phi(x) - \tilde{\Phi}(x) \right| = \sup_{x \in [-1,1]^d} \left| A_2 \sigma(A_1 x) - \tilde{A}_2 \sigma(\tilde{A}_1 x) \right|$$
$$\leq \frac{1}{4} N d \, 2^{-b} + N 2^{-b}.$$

Problem 4 (25 points)

- (a) Let X_1 be a finite subset of \mathbb{R}^d , $d \in \mathbb{N}$, let $\{X_1^+, X_1^-\}$ be a dichotomy of X_1 , and let $\phi : \mathbb{R}^d \mapsto \mathbb{R}^m$, $m \in \mathbb{N}$, be a mapping. State the definition for the dichotomy $\{X_1^+, X_1^-\}$ to be homogeneously linearly separable and the definition for it to be ϕ -separable. Show that there are at most $2^{\operatorname{card}(X_1)}$ homogeneously linearly separable dichotomies of X_1 , with equality if and only if every dichotomy is homogeneously linearly separable. Here, $\operatorname{card}(X_1)$ denotes the cardinality of X_1 .
- (b) Consider a vector $x_1 \in \mathbb{R}^d$ and let *C* be the number of homogeneously linearly separable dichotomies of X_1 . Show that the number of homogeneously linearly separable dichotomies of $X_1 \cup \{x_1\}$ is at most 2*C*.
- (c) Consider the set $X_2 = \{(1,0), (-1,0), (0,1)\} \subset \mathbb{R}^2$ and its dichotomy

$$\{X_2^+ = \{(1,0), (-1,0)\}, X_2^- = \{(0,1)\}\}.$$

Is this dichotomy homogeneously linearly separable? Justify your answer.

- (d) Consider the mapping $\phi_1 : \mathbb{R}^2 \mapsto \mathbb{R}^2$ given by $\phi_1(x, y) = (x^2, y^2)$. Is the dichotomy $\{X_2^+ = \{(1, 0), (-1, 0)\}, X_2^- = \{(0, 1)\}\}$ from subproblem (c) ϕ_1 -separable? Justify your answer.
- (e) Consider $X_3 = \{(1,0), (0,1), (-1,0), (0,-1)\}$, with $X_3^+ = \{(-1,0), (1,0)\}$ and $X_3^- = \{(0,1), (0,-1)\}$. Let $\phi_2 : \mathbb{R}^2 \mapsto \mathbb{R}^4$ be given by $\phi_2(x,y) = (x^2 + y^2, x, y, 1)$. Show that the dichotomy $\{X_3^+, X_3^-\}$ is not ϕ_2 -separable.