**ETH**

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

**Prof. Dr. Helmut Bölcskei**
**Chair for Mathematical Information Science**

# Solutions to the
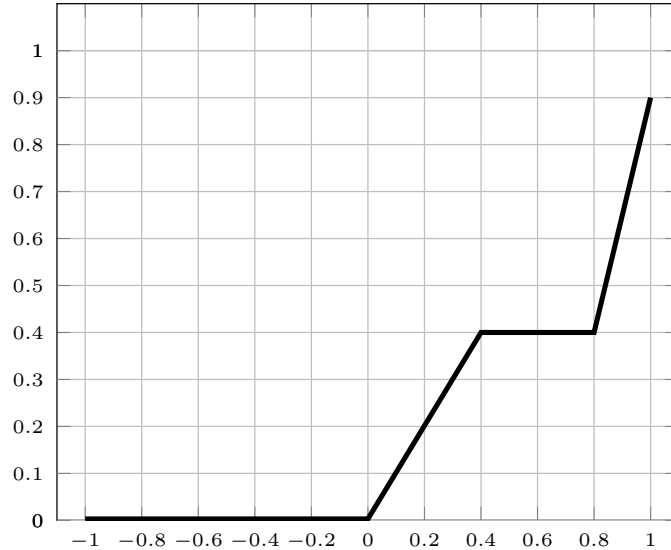# Exam on Neural Network Theory
# February 11, 2021

## Problem 1

(a) Note that $|x| = \rho(x) + \rho(-x)$ and $10x = 1.25\rho \circ 2\rho \circ 2\rho \circ 2\rho(x), \forall x \geq 0$. Since $|x| \geq 0$, the function $f(x) = 10|x|$ can hence be realized through the network

$$\Psi(x) = 1.25 \circ \rho \circ 2 \circ \rho \circ 2 \circ \rho \circ \begin{pmatrix} 2 & 2 \end{pmatrix} \circ \rho \circ \begin{pmatrix} 1 \\ -1 \end{pmatrix} x.$$

This network satisfies $\mathcal{W}(\Psi) = 2$, $\mathcal{B}(\Psi) = 2$, and has depth $\mathcal{L}(\Psi) = 5$ and connectivity $\mathcal{M}(\Psi) = 7$.
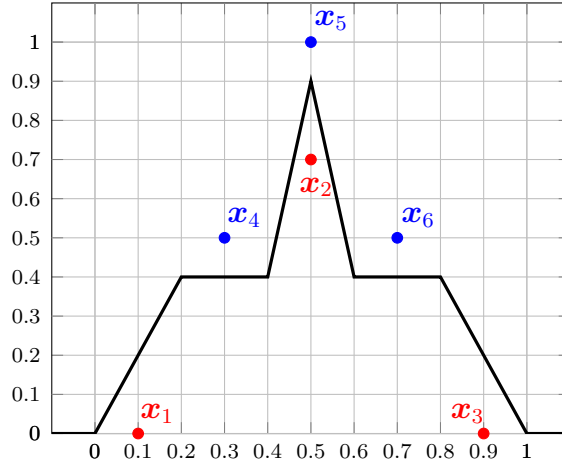
(b) The network $\Phi$ consists of two affine maps, its depth is therefore $2$, the width is $3$ and the connectivity is $8$. The plot of $\Phi(x)$ on the interval $[-1, 1]$ is depicted below.



(c) The function $f(x)$ can be inferred from its definition according to $f(x) = \Phi(2x)$ on $[-\infty, 1/2]$ and by exploiting symmetry. In summary, we get

$$f(x) = \begin{cases} \Phi(2x), & \text{if } x \leq 1/2, \\ \Phi(2 - 2x), & \text{if } x > 1/2, \end{cases}$$

where $\Phi(x)$ is the ReLU network from subproblem (b). The following figure depicting $f(x)$ shows that $f(x)$ is class-separating for the problem at hand.

Denote the ReLU function as $\rho(x) = \max\{0, x\}$. The function $f(x)$ can be realized according to

$$f(x) = \rho(2x) - \rho(2x - 0.4) + \rho(5x - 2) - \rho(10x - 5)$$
$$+ \rho(5x - 3) - \rho(2x - 1.6) + \rho(2x - 2).$$

The corresponding depth-2 ReLU network is given by

$$\Phi(x) = \big(W_2 \circ \rho \circ W_1\big)(x)$$

with

$$W_1(x) = \begin{pmatrix} 2 \\ 2 \\ 5 \\ 10 \\ 5 \\ 2 \\ 2 \end{pmatrix} x + \begin{pmatrix} 0 \\ -0.4 \\ -2 \\ -5 \\ -3 \\ -1.6 \\ -2 \end{pmatrix}, \quad W_2(x) = \begin{pmatrix} 1 & -1 & 1 & -1 & 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \end{pmatrix}.$$

The network $\Phi$ has depth 2, width 7, and connectivity 20.

(d) The sawtooth functions $g_s(x)$ are periodic with period $2^{-s+1}$, see Figure 1. Thus, the cardinality of the set $\{x : g_s(x) = 1, \ x \in [0, 1]\}$ is $2^{s-1}$.
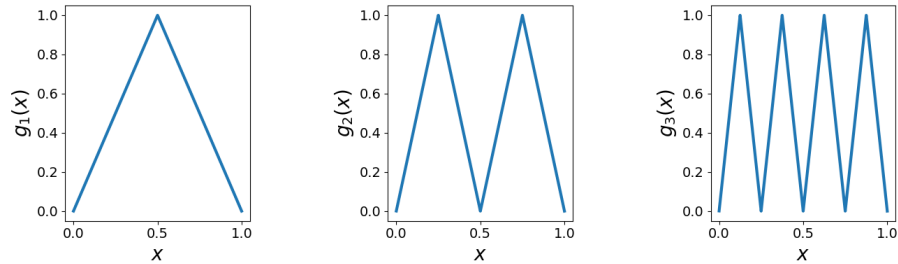


Figure 1: Sawtooth functions

(e) Note that

$$g(x) = \rho(2x) - \rho(4x - 2) + \rho(2x - 2) = \begin{cases} 2x, & \text{if } 0 \leq x \leq 1/2, \\ 2 - 2x, & \text{if } 1/2 < x \leq 1, \\ 0, & \text{otherwise.} \end{cases}$$

It therefore follows that

$$f(x) = \big(\Phi \circ g\big)(x),$$

where we also used that $\Phi(x) = 0, \forall x \leq 0$. The sawtooth function $g(x)$ can be realized through a ReLU network according to $g(x) = \rho(2x) - \rho(4x - 2) + \rho(2x - 2) = \big(W_2^g \circ \rho \circ W_1^g\big)(x)$ with

$$W_1^g(x) = \begin{pmatrix} 2 \\ 4 \\ 2 \end{pmatrix} x + \begin{pmatrix} 0 \\ -2 \\ -2 \end{pmatrix}, \quad W_2^g(x) = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$

Now, recall that the network $\Phi(x)$ is given by $\Phi(x) = \big(W_2 \circ \rho \circ W_1\big)(x)$ with

$$W_1(x) = \begin{pmatrix} 1 \\ 1 \\ 2.5 \end{pmatrix} x - \begin{pmatrix} 0 \\ 0.4 \\ 2 \end{pmatrix}, \quad W_2(x) = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$

The composition $f(x) = \big(\Phi \circ g\big)(x)$ can hence be realized through the network

$$\Phi_2(x) = \big(W_2 \circ \rho \circ W_1 W_2^g \circ \rho \circ W_1^g\big)(x) = \big(W_3' \circ \rho \circ W_2' \circ \rho \circ W_1'\big)(x),$$

where

$$W_1'(x) = W_1^g(x) = \begin{pmatrix} 2 \\ 4 \\ 2 \end{pmatrix} x + \begin{pmatrix} 0 \\ -2 \\ -2 \end{pmatrix}, \quad W_3'(x) = W_2(x) = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix},$$

$$W_2'(x) = \big(W_1 W_2^g\big)(x) = \begin{pmatrix} 1 \\ 1 \\ 2.5 \end{pmatrix} \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} - \begin{pmatrix} 0 \\ 0.4 \\ 2 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & -1 & 1 \\ 1 & -1 & 1 \\ 2.5 & -2.5 & 2.5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} - \begin{pmatrix} 0 \\ 0.4 \\ 2 \end{pmatrix}.$$

The network $\Phi_2$ has depth 3, width 3, and connectivity 19.

# Problem 2

(a) (i) An $\epsilon$-covering of the compact set $\mathcal{C} \subseteq \mathcal{X}$ with respect to the metric $\rho$ is a set $\{x_1, \ldots, x_N\} \subset \mathcal{C}$ such that for each $x \in \mathcal{C}$, there exists an $i \in [1, N]$ so that $\rho(x, x_i) \leq \epsilon$. The $\epsilon$-covering number $N(\epsilon; \mathcal{C}, \rho)$ is the cardinality of a smallest $\epsilon$-covering of $\mathcal{C}$.

(ii) An $\epsilon$-packing of a compact set $\mathcal{C} \subseteq \mathcal{X}$ with respect to the metric $\rho$ is a set $\{x_1, \ldots, x_N\} \subset \mathcal{C}$ such that $\rho(x_i, x_j) > \epsilon$, for all distinct $i, j$. The $\epsilon$-packing number $M(\epsilon; \mathcal{C}, \rho)$ is the cardinality of a largest $\epsilon$-packing of $\mathcal{C}$.

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

(b) (i) For every $f_\theta \in \mathcal{F}$, we can find a $\theta_i$ in the set $\{\theta_0, \ldots, \theta_T, \theta_{T+1}\}$, such that $|\theta_i - \theta| \leq \epsilon$. We then have

$$\|f_{\theta_i} - f_\theta\|_\infty = \max_{x \in [0,1]} \left| \ln(1 + \theta_i x) - \ln(1 + \theta x) \right| = \max_{x \in [0,1]} \left| \ln\left(\frac{1 + \theta_i x}{1 + \theta x}\right) \right|$$

$$= \max_{x \in [0,1]} \left| \ln\left(1 + \frac{(\theta_i - \theta)x}{1 + \theta x}\right) \right| \leq \max_{x \in [0,1]} \left| \frac{(\theta_i - \theta)x}{1 + \theta x} \right| \leq |\theta_i - \theta| \leq \epsilon.$$

Therefore, we can conclude that the set $\{f_{\theta_0}, \ldots, f_{\theta_T}, f_{\theta_{T+1}}\}$ constitutes an $\epsilon$-covering of $\mathcal{F}$. An upper bound on the covering number is hence given by $N(\epsilon; \mathcal{F}, \|\cdot\|_\infty) \leq T + 2 \leq \frac{1}{2\epsilon} + 2$.

(ii) We construct an explicit packing as follows. Set $T = \lfloor \frac{1}{3\epsilon} \rfloor$, and for $i = 0, 1, \ldots, T$, define $\theta_i = 3\epsilon i$. Moreover, note that for all $i, j$ with $i \neq j$, we have

$$\|f_{\theta_i} - f_{\theta_j}\|_\infty = \max_{x \in [0,1]} \left| \ln(1 + \theta_i x) - \ln(1 + \theta_j x) \right| = \max_{x \in [0,1]} \left| \ln\left(\frac{1 + \theta_i x}{1 + \theta_j x}\right) \right|$$

$$= \max_{x \in [0,1]} \left| \ln\left(1 + \frac{(\theta_i - \theta_j)x}{1 + \theta_j x}\right) \right| \geq \max_{x \in [0,1]} \left| \left(\frac{(\theta_i - \theta_j)x}{1 + \theta_j x}\right) \middle/ \left(1 + \frac{(\theta_i - \theta_j)x}{1 + \theta_j x}\right) \right|$$

$$= \max_{x \in [0,1]} \left| \left(\frac{(\theta_i - \theta_j)x}{1 + \theta_j x}\right) \middle/ \left(\frac{1 + \theta_j x + \theta_i x - \theta_j x}{1 + \theta_j x}\right) \right| = \max_{x \in [0,1]} \left| \frac{(\theta_i - \theta_j)x}{1 + \theta_i x} \right|$$

$$\geq \max_{x \in [0,1]} \left| \frac{(\theta_i - \theta_j)x}{2} \right| = \left| \frac{\theta_i - \theta_j}{2} \right| = \left| \frac{3\epsilon(i - j)}{2} \right| > \epsilon,$$

by definition of $\theta_i$. We can therefore conclude that $\{f_{\theta_0}, \ldots, f_{\theta_T}\}$ is an $\epsilon$-packing and the corresponding packing number satisfies $M(\epsilon; \mathcal{F}, \|\cdot\|_\infty) \geq T + 1 \geq \frac{1}{3\epsilon}$.

(iii) By subproblems (a.iii), (b.i), and (b.ii), we obtain

$$\frac{1}{6\epsilon} \leq M(2\epsilon; \mathcal{F}, \|\cdot\|_\infty) \leq N(\epsilon; \mathcal{F}, \|\cdot\|_\infty) \leq \frac{1}{2\epsilon} + 2,$$

which allows us to conclude that $\log N(\epsilon; \mathcal{F}, \|\cdot\|_\infty) \asymp \log(1/\epsilon)$, as $\epsilon \to 0$.

4

# Problem 3

(a) As $a$ and $T_b(a)$ differ only in digits after the $b$-th position in the fractional parts of their binary representation, we have

$$|a - \tilde{a}| \le \sum_{i=b+1}^{\infty} 2^{-i} = 2^{-b} \sum_{i=1}^{\infty} 2^{-i} = 2^{-b},$$

where we used $\sum_{i=1}^{\infty} 2^{-i} = 1$ in the last equality.

(b) As $T_b$ acts entrywise on $A_1$, we have

$$\left\| A_1 - \tilde{A}_1 \right\|_{\infty} = \| A_1 - T_b(A_1) \|_{\infty} \le \sup_{a \in \mathbb{R}} |a - T_b(a)| \le 2^{-b}. \tag{1}$$

We therefore get

$$\sup_{x \in [-1,1]^d} \left\| Ax - \tilde{A}_1 x \right\|_{\infty} = \sup_{x \in [-1,1]^d} \left\| \left( A_1 - \tilde{A}_1 \right) x \right\|_{\infty} \tag{2}$$

$$\le \sup_{x \in [-1,1]^d} d \left\| A_1 - \tilde{A}_1 \right\|_{\infty} \| x \|_{\infty} \tag{3}$$

$$\le d \, 2^{-b}, \tag{4}$$

where in (3) we used (1) and the inequality in the hint.

(c) The first derivative of $\sigma$ is given by $\sigma'(x) = \frac{e^x}{(1+e^x)^2}$, $x \in \mathbb{R}$, which is positive for all $x \in \mathbb{R}$. It therefore follows that $\sigma$ is strictly increasing, and hence $\sup_{x \in \mathbb{R}} \sigma(x) \le \lim_{x \to \infty} \sigma(x) = 1$ and $\inf_{x \in \mathbb{R}} \sigma(x) \ge \lim_{x \to -\infty} \sigma(x) = 0$. The second derivative of $\sigma$ is given by $\sigma''(x) = \frac{e^x(1-e^x)}{(1+e^x)^3}$, $x \in \mathbb{R}$, which is positive on $(-\infty, 0)$ and negative on $(0, \infty)$. This implies that $\sigma'$ is increasing on $(-\infty, 0)$ and decreasing on $(0, \infty)$. Hence $\sigma'$ attains its maximum at $x = 0$ with $\sigma'(0) = \frac{1}{4}$. For every $x, y \in \mathbb{R}$, by the mean value theorem, there exists a real number $z$ between $x$ and $y$ such that

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

which implies

$$|\sigma(x) - \sigma(y)| = |\sigma'(z)(x - y)| \le \frac{1}{4} |x - y|,$$

where we used $0 < \sigma'(z) \le \sigma'(0) = \frac{1}{4}$.

(d) We have

$$\sup_{x \in [-1,1]^d} \left\| \sigma \left( A_1 x \right) - \sigma \left( \tilde{A}_1 x \right) \right\|_{\infty} \tag{5}$$

$$\le \sup_{x \in [-1,1]^d} \frac{1}{4} \left\| A_1 x - \tilde{A}_1 x \right\|_{\infty} \tag{6}$$

$$\le \frac{1}{4} d \, 2^{-b}, \tag{7}$$

where in (6) we used the definition of the $\|\cdot\|_{\infty}$-norm and the Lipschitz continuity established in subproblem (c), and (7) follows from subproblem (b).

5

(e) We have, for all $x \in [-1, 1]^d$,

$$\left| \Phi(x) - \tilde{\Phi}(x) \right| \tag{8}$$

$$= \left\| A_2 \sigma(A_1 x) - \tilde{A}_2 \sigma(\tilde{A}_1 x) \right\|_\infty \tag{9}$$

$$= \left\| A_2 \sigma(A_1 x) - A_2 \sigma(\tilde{A}_1 x) + A_2 \sigma(\tilde{A}_1 x) - \tilde{A}_2 \sigma(\tilde{A}_1 x) \right\|_\infty \tag{10}$$

$$\leq \left\| A_2 \left( \sigma(A_1 x) - \sigma(\tilde{A}_1 x) \right) \right\|_\infty + \left\| (A_2 - \tilde{A}_2) \sigma(\tilde{A}_1 x) \right\|_\infty \tag{11}$$

$$\leq N \left\| A_2 \right\|_\infty \left\| \sigma(A_1 x) - \sigma(\tilde{A}_1 x) \right\|_\infty + N \left\| A_2 - \tilde{A}_2 \right\|_\infty \left\| \sigma(\tilde{A}_1 x) \right\|_\infty \tag{12}$$

$$\leq \frac{1}{4} N d\, 2^{-b} + N 2^{-b}, \tag{13}$$

where in (12) we used the inequality derived in (3), and in (13) we employed $\|A_2\|_\infty \leq 1$, the result from subproblem (d), and $\left\| \sigma(\tilde{A}_1 x) \right\|_\infty \leq 1$.

# Problem 4

(a) The dichotomy $\{X_1^+, X_1^-\}$ is said to be homogeneously linearly separable if there exists a nonzero vector $w_1 \in \mathbb{R}^d$ such that

$$\langle x, w_1 \rangle > 0, \text{ for all } x \in X_1^+,$$
$$\langle x, w_1 \rangle < 0, \text{ for all } x \in X_1^-,$$

and it is said to be $\phi$-separable if there exists a nonzero vector $w_2 \in \mathbb{R}^m$ such that

$$\langle \phi(x), w_2 \rangle > 0, \text{ for all } x \in X_1^+,$$
$$\langle \phi(x), w_2 \rangle < 0, \text{ for all } x \in X_1^-.$$

Since there are at most $2^{\mathrm{card}(X_1)}$ dichotomies of $X_1$, as explained in the lecture, it follows from the inclusion relation that the number of homogeneously linearly separable dichotomies of $X_1$ is less than or equal to $2^{\mathrm{card}(X_1)}$, with equality if every dichotomy of $X_1$ is homogeneously linearly separable.

(b) Let $S_{X_1}$, $S_{X_1 \cup \{x_1\}}$ be the sets of homogeneously linearly separable dichotomies of $X_1$ and $X_1 \cup \{x_1\}$, respectively. Consider a dichotomy $\{X^+, X^-\} \in S_{X_1 \cup \{x_1\}}$ and note that $\{X^+, X^-\}$ can be written as

$$\left\{X_1^+ \cup \{x_1\}, X_1^-\right\} \text{ or } \left\{X_1^+, X_1^- \cup \{x_1\}\right\} \tag{14}$$

for some dichotomy $\{X_1^+, X_1^-\}$ of $X_1$. As $X_1^+ \subset X^+$ and $X_1^- \subset X^-$, the dichotomy $\{X_1^+, X_1^-\}$ can be separated by the hyperplane that separates $\{X^+, X^-\}$. Therefore, the dichotomy $\{X_1^+, X_1^-\}$ in (14) is homogeneously linearly separable and we have

$$S_{X_1 \cup \{x_1\}} \subset \left\{\left\{X_1^+ \cup \{x_1\}, X_1^-\right\} \text{ or } \left\{X_1^+, X_1^- \cup \{x_1\}\right\} : \left\{X_1^+, X_1^-\right\} \in S_{X_1}\right\},$$

which implies

$$\mathrm{card}\left(S_{X_1 \cup \{x_1\}}\right) \leq 2 \, \mathrm{card}\left(S_{X_1}\right) = 2C.$$

(c) The dichotomy is not homogeneously linearly separable as there is no line through the origin such that $(-1, 0)$ and $(1, 0)$ lie on the same side of the line. Otherwise, there would exist a nonzero vector $w = (w_1, w_2) \in \mathbb{R}^2$ such that $\langle w, (1, 0) \rangle = w_1 > 0$ and $\langle w, (-1, 0) \rangle = -w_1 < 0$, which constitutes a contradiction. See Fig. 2 for an illustration.
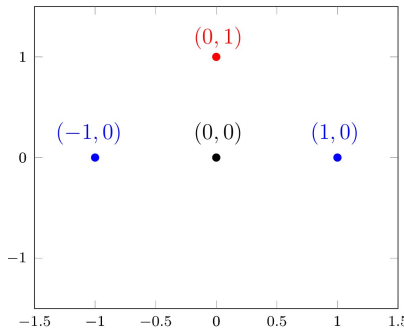


Figure 2: If the dichotomy were homogeneously linearly separable, there would exist a line through the origin that separates the blue points from the red point.

(d) The dichotomy is $\phi_1$-separable. Let $w = (1, -1)$. Then, $\langle \phi_1(1, 0), w \rangle = 1 > 0$, $\langle \phi_1(-1, 0), w \rangle = 1 > 0$, and $\langle \phi_1(0, 1), w \rangle = -1 < 0$.

(e) Suppose for the sake of contradiction that the dichotomy $\{ X_3^+, X_3^- \}$ is $\phi_2$-separable. Then, there would exist a nonzero vector $w = (w_1, w_2, w_3, w_4) \in \mathbb{R}^4$ such that

$$\langle \phi_2(x, y), w \rangle > 0, \quad \text{for all } (x, y) \in X_3^+,$$
$$\langle \phi_2(x, y), w \rangle < 0, \quad \text{for all } (x, y) \in X_3^-,$$

that is

$$
\begin{aligned}
w_1 + w_2 \qquad + w_4 &> 0, & (15) \\
w_1 - w_2 \qquad + w_4 &> 0, & (16) \\
w_1 \qquad + w_3 + w_4 &< 0, & (17) \\
w_1 \qquad - w_3 + w_4 &< 0. & (18)
\end{aligned}
$$

Adding (15) and (16), we obtain $w_1 + w_4 > 0$, while adding (17) and (18) yields $w_1 + w_4 < 0$, which establishes the contradiction.