

# Solutions to the Exam on Neural Network Theory February 6, 2024

## Problem 1

(a) One possible solution is

$$\begin{aligned}\Phi(x) &= \rho(x) - \rho(x-1) \\ &= (1 \quad -1) \circ \rho \circ \left( \begin{pmatrix} 1 \\ 1 \end{pmatrix} x + \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right), \quad x \in \mathbb{R}.\end{aligned}$$

This network satisfies  $\mathcal{L}(\Phi) = 2$ ,  $\mathcal{W}(\Phi) = 2$ , and  $\mathcal{M}(\Phi) = 5$ .

(b) We directly calculate the function  $f$  as

$$f(x) = \begin{cases} 0, & \text{for } x < 0 \\ 4x, & \text{for } 0 \leq x < \frac{1}{4} \\ -2x + \frac{3}{2}, & \text{for } \frac{1}{4} \leq x < \frac{3}{4} \\ 4x - 3, & \text{for } \frac{3}{4} \leq x < 1 \\ 1, & \text{for } x \geq 1. \end{cases}$$

The sketch of  $f$  is given below.

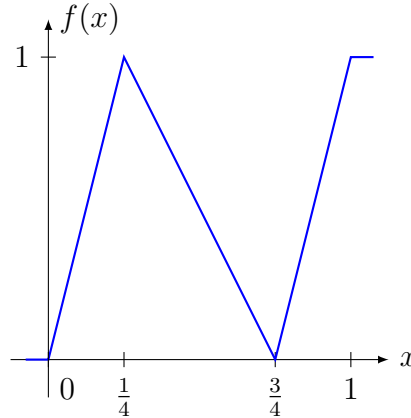


Figure 1:  $f(x)$ .

One possible solution for  $\Phi^f$  is

$$\begin{aligned}\Phi^f(x) &= 4\rho(x) - 6\rho(x-1/4) + 6\rho(x-3/4) - 4\rho(x-1) \\ &= (4 \quad -6 \quad 6 \quad -4) \circ \rho \circ \left( \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} x + \begin{pmatrix} 0 \\ -1/4 \\ -3/4 \\ -1 \end{pmatrix} \right), \quad x \in \mathbb{R},\end{aligned}$$

which satisfies  $\mathcal{L}(\Phi^f) = 2$ ,  $\mathcal{W}(\Phi^f) = 4$ ,  $\mathcal{M}(\Phi^f) = 11$ , and  $\mathcal{B}(\Phi^f) = 6$ .

(c) We first realize the function  $g$  as a  $\sigma$ -network according to

$$g(x, y) = (W_3 \circ \sigma \circ W_2 \circ \sigma \circ W_1)(x, y), \quad (1)$$

where

$$\begin{aligned} W_1(x, y) &= \begin{pmatrix} -2 & -1 \\ 0.5 & -2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad x, y \in \mathbb{R}, \\ W_2(x) &= \begin{pmatrix} 1 & 1 \end{pmatrix} x, \quad x \in \mathbb{R}^2, \\ W_3(x) &= x, \quad x \in \mathbb{R}. \end{aligned}$$

Based on the result from subproblem (a), namely

$$\sigma(x) = \rho(x) - \rho(x - 1),$$

we can convert  $\sigma$  in (1) to obtain a ReLU network realizing  $g(x, y)$  according to

$$\Phi^g(x, y) = (W'_3 \circ \rho \circ W'_2 \circ \rho \circ W'_1)(x, y), \quad (2)$$

where

$$\begin{aligned} W'_1(x, y) &= \begin{pmatrix} -2 & -1 \\ -2 & -1 \\ 0.5 & -2 \\ 0.5 & -2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ 0 \\ -1 \end{pmatrix}, \\ W'_2(x) &= \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & -1 & 1 & -1 \end{pmatrix} x + \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \quad x \in \mathbb{R}^4, \\ W'_3(x) &= \begin{pmatrix} 1 & -1 \end{pmatrix} x, \quad x \in \mathbb{R}^2. \end{aligned}$$

Inspection of (2) shows that  $\mathcal{L}(\Phi^g) = 3$  as desired.

(d) We can express the operation  $\vee$  by affine copies of  $\rho$  according to

$$x \vee y = \max\{x, y\} = x + \max\{0, y - x\} = \rho(x) - \rho(-x) + \rho(y - x), \quad \text{for } x, y \in \mathbb{R}.$$

One possible solution for  $\Phi^\vee$  is hence

$$\Phi^\vee(x, y) = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} \circ \rho \circ \left( \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \right), \quad x, y \in \mathbb{R}.$$

Likewise, for the operation  $\wedge$ , we have

$$x \wedge y = \min\{x, y\} = x - \max\{0, x - y\} = \rho(x) - \rho(-x) - \rho(x - y), \quad \text{for } x, y \in \mathbb{R}.$$

So we can choose  $\Phi^\wedge$  according to

$$\Phi^\wedge(x, y) = \begin{pmatrix} 1 & -1 & -1 \end{pmatrix} \circ \rho \circ \left( \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \right), \quad x, y \in \mathbb{R}.$$

(e) In (d), the "min" operation can be realized according to

$$\min\{x, y\} = (W_2 \circ \rho \circ W_1)(x, y),$$

where

$$W_1(x, y) = \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}, \quad x, y \in \mathbb{R},$$

$$W_2(x) = (1 \quad -1 \quad -1) x, \quad x \in \mathbb{R}^3.$$

Using  $z = \rho(z) - \rho(-z)$ , for  $z \in \mathbb{R}$ , the ReLU network  $W_4 \circ \rho \circ W_3$  with

$$W_3(x, y, z) = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & -1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}, \quad x, y, z \in \mathbb{R},$$

$$W_4(x) = \begin{pmatrix} 1 & -1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \end{pmatrix} x, \quad x \in \mathbb{R}^5.$$

satisfies

$$(W_4 \circ \rho \circ W_3)(x, y, z) = \begin{pmatrix} \min\{x, y\} \\ z \end{pmatrix}, \quad x, y, z \in \mathbb{R}.$$

Let  $\Phi^h := W_2 \circ \rho \circ W_1 \circ W_4 \circ \rho \circ W_3$ . Then  $\Phi^h(x, y, z) = \min\{\min\{x, y\}, z\} = \min\{x, y, z\}$  as required.

## Problem 2

- (a) An  $\epsilon$ -covering of  $\mathcal{F}$  with respect to the metric  $\rho_\infty$  is a set  $\{x_1, \dots, x_N\} \subset \mathcal{F}$  such that for each  $x \in \mathcal{F}$ , there exists an  $i \in \{1, \dots, N\}$  so that  $\rho_\infty(x, x_i) \leq \epsilon$ . The  $\epsilon$ -covering number  $N(\epsilon; \mathcal{F}, \rho_\infty)$  is the cardinality of a smallest  $\epsilon$ -covering of  $\mathcal{F}$ .

- (b) For all  $\theta, \theta' \in [0, 1]$ , we have

$$\rho_\infty(f_{\theta, \theta'}, f_{0,0}) = \sup_{x \in [0,1]} |f_{\theta, \theta'}(x) - f_{0,0}(x)| = \max_{x \in [0,1]} |1 - e^{-\theta x} + \theta'| = \max_{x \in [0,1]} (1 - e^{-\theta x} + \theta') \leq 2.$$

Therefore, for  $\epsilon \geq 2$ , the  $\epsilon$ -ball around  $f_{0,0}$  contains all elements in  $\mathcal{F}$ . The singleton set  $\{f_{0,0}\}$  constitutes an  $\epsilon$ -covering of  $\mathcal{F}$ , which establishes

$$N(\epsilon; \mathcal{F}, \rho_\infty) = 1, \quad \text{for } \epsilon \geq 2.$$

- (c) For given  $\epsilon < 2$ , for every  $f_{\theta, \theta'} \in \mathcal{F}$ , we can find  $(\theta_i, \theta'_j)$  in the set  $\{(\theta_i, \theta'_j) : i, j = 0, 1, \dots, T+1\}$ , such that  $|\theta_i - \theta| \leq \epsilon/2$  and  $|\theta'_j - \theta'| \leq \epsilon/2$ . We then have

$$\begin{aligned} \rho(f_{\theta, \theta'} - f_{\theta_i, \theta'_j}) &= \sup_{x \in [0,1]} |f_{\theta, \theta'}(x) - f_{\theta_i, \theta'_j}(x)| \\ &= \max_{x \in [0,1]} |-e^{-\theta x} + e^{-\theta_i x} + \theta' - \theta'_j| \\ &\leq |\theta' - \theta'_j| + \max_{x \in [0,1]} |e^{-\theta x} - e^{-\theta_i x}| \\ &\leq |\theta' - \theta'_j| + |\theta_i - \theta| \leq \epsilon, \end{aligned}$$

where we used, for  $\theta < \theta_i$ ,

$$\begin{aligned} \max_{x \in [0,1]} |e^{-\theta x} - e^{-\theta_i x}| &= \max_{x \in [0,1]} e^{-\theta x} |1 - e^{-(\theta_i - \theta)x}| \\ &\leq \max_{x \in [0,1]} (1 - e^{-(\theta_i - \theta)x}) \\ &\leq \max_{x \in [0,1]} (\theta_i - \theta)x \\ &= |\theta_i - \theta|. \end{aligned}$$

The case  $\theta > \theta_i$  follows similarly. We conclude that the set  $\{f_{\theta_i, \theta'_j} : i, j = 0, \dots, T+1\}$  constitutes an  $\epsilon$ -covering of  $\mathcal{F}$ . An upper bound on the covering number is hence given by  $N(\epsilon; \mathcal{F}, \rho_\infty) \leq (T+2)^2 \leq (\frac{1}{\epsilon} + 2)^2$ .

- (d) We construct an  $\epsilon$ -packing as follows. Set  $\theta_0 = 0$  and define  $\theta_i = -\log(1 - \epsilon i)$  for all  $i$  such that  $\theta_i \leq 1$ . The largest index  $T$  so that this holds is given by  $T = \left\lfloor \frac{1 - e^{-1}}{\epsilon} \right\rfloor$ . For  $j = 0, 1, \dots, \left\lfloor \frac{1}{\epsilon} \right\rfloor$ , let  $\theta'_j = j\epsilon$ . Note that for any two distinct points  $(\theta_i, \theta'_j)$  and  $(\theta_m, \theta'_n)$ , if  $j \neq n$ , we have

$$\rho(f_{\theta_i, \theta'_j}, f_{\theta_m, \theta'_n}) = \max_{x \in [0,1]} |f_{\theta_i, \theta'_j}(x) - f_{\theta_m, \theta'_n}(x)| \geq |f_{\theta_i, \theta'_j}(0) - f_{\theta_m, \theta'_n}(0)| = |\theta'_j - \theta'_n| \geq \epsilon,$$

and if  $j = n$ , then  $i \neq m$ , and we have

$$\rho(f_{\theta_i, \theta'_j}, f_{\theta_m, \theta'_n}) = \max_{x \in [0,1]} |f_{\theta_i, \theta'_j}(x) - f_{\theta_m, \theta'_n}(x)| \geq |f_{\theta_i, \theta'_j}(1) - f_{\theta_m, \theta'_n}(1)| = |\epsilon(i - m)| \geq \epsilon.$$

We can therefore conclude that the set  $\{f_{\theta_i, \theta'_j} : i = 0, \dots, T, j = 0, \dots, \left\lfloor \frac{1}{\epsilon} \right\rfloor\}$  is an  $\epsilon$ -packing and the packing number satisfies  $M(\epsilon; \mathcal{F}, \rho_\infty) \geq (T+1)(1 + \left\lfloor \frac{1}{\epsilon} \right\rfloor) \geq \frac{1 - e^{-1}}{\epsilon^2}$ .

(e) By subproblems (c) and (d) and Lemma 1 in the Handout, we obtain

$$\frac{1 - e^{-1}}{(2\epsilon)^2} \leq M(2\epsilon; \mathcal{F}, \rho_\infty) \leq N(\epsilon; \mathcal{F}, \rho_\infty) \leq \left(\frac{1}{\epsilon} + 2\right)^2,$$

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

### Problem 3

(a)

$$\text{Softmax}^{(1)}(x) = \frac{\exp(x)}{\exp(x)} = 1, \quad (3)$$

for all  $x \in \mathbb{R}$ .

(b) Fix  $x \in \mathbb{R}^n$  and  $i \in \{1, \dots, n\}$ . For all  $t \in \mathbb{R}$ ,  $\exp(t) \geq 0$ , so

$$\text{Softmax}_i^{(n)}(x) = \frac{\exp(x_i)}{\sum_{k=1}^n \exp(x_k)} \geq 0. \quad (4)$$

Moreover,

$$\text{Softmax}_i^{(n)}(x) = \frac{\exp(x_i)}{\sum_{k=1}^n \exp(x_k)} \leq \frac{\exp(x_i)}{\exp(x_i)} \leq 1. \quad (5)$$

(c) Let  $n \in \mathbb{N}$ ,  $i, j \in \{1, \dots, n\}$ , and  $x \in \mathbb{R}^n$ . Suppose that  $i \neq j$ .

$$\partial_j \text{Softmax}_i^{(n)}(x) = \partial_j \left( \frac{\exp(x_i)}{\sum_{k=1}^n \exp(x_k)} \right) \quad (6)$$

$$= - \frac{\exp(x_i) \partial_j (\sum_{k=1}^n \exp(x_k))}{(\sum_{k=1}^n \exp(x_k))^2} \quad (7)$$

$$= - \frac{\exp(x_i) \exp(x_j)}{(\sum_{k=1}^n \exp(x_k))^2} \quad (8)$$

$$= - \text{Softmax}_i^{(n)}(x) \text{Softmax}_j^{(n)}(x). \quad (9)$$

$$\partial_i \text{Softmax}_i^{(n)}(x)_i = \partial_i \left( \frac{\exp(x_i)}{\sum_{k=1}^n \exp(x_k)} \right) \quad (10)$$

$$= \frac{(\sum_{k=1}^n \exp(x_k)) \partial_i (\exp(x_i)) - \exp(x_i) \partial_i (\sum_{k=1}^n \exp(x_k))}{(\sum_{k=1}^n \exp(x_k))^2} \quad (11)$$

$$= \frac{(\sum_{k=1}^n \exp(x_k)) \exp(x_i) - \exp(x_i)^2}{(\sum_{k=1}^n \exp(x_k))^2} \quad (12)$$

$$= \frac{\exp(x_i)}{\sum_{k=1}^n \exp(x_k)} - \frac{\exp(x_i)^2}{(\sum_{k=1}^n \exp(x_k))^2} \quad (13)$$

$$= \text{Softmax}_i^{(n)}(x) - \text{Softmax}_i^{(n)}(x) \text{Softmax}_i^{(n)}(x). \quad (14)$$

Combining (14) for each  $i \in \{1, \dots, n\}$ , we obtain the result in the hint, namely

$$\nabla \text{Softmax}^{(n)}(x) = \text{diag} \left( \text{Softmax}^{(n)}(x) \right) - \text{Softmax}^{(n)}(x) \text{Softmax}^{(n)}(x)^T. \quad (15)$$

By subproblem b), one has  $0 \leq \text{Softmax}_i^{(n)}(x) \leq 1$ , for all  $x \in \mathbb{R}^n$ ,  $i \in \{1, \dots, n\}$ . It follows that  $-1 \leq \nabla \text{Softmax}_{i,j}^{(n)}(x) \leq 1$ , for all  $x \in \mathbb{R}^n$ ,  $i, j \in \{1, \dots, n\}$ . Therefore,  $\|\nabla \text{Softmax}^{(n)}(x)\|_\infty \leq 1$ , for all  $x \in \mathbb{R}^n$ , as desired.

(d) Let  $n \in \mathbb{N}$  and  $x, y \in \mathbb{R}^n$ ,  $x \neq y$ . By Theorem 1 in the Handout, there exists  $z \in \mathbb{R}^n$  such that

$$\|\text{Softmax}^{(n)}(y) - \text{Softmax}^{(n)}(x)\|_2 \leq \|\nabla \text{Softmax}^{(n)}(z)(y - x)\|_2. \quad (16)$$

It follows that

$$\|\text{Softmax}^{(n)}(y) - \text{Softmax}^{(n)}(x)\|_\infty \stackrel{(a)}{\leq} \|\text{Softmax}^{(n)}(y) - \text{Softmax}^{(n)}(x)\|_2 \quad (17)$$

$$\stackrel{(16)}{\leq} \|\nabla \text{Softmax}^{(n)}(z)(y - x)\|_2 \quad (18)$$

$$\stackrel{(b)}{\leq} \sqrt{n} \|\nabla \text{Softmax}^{(n)}(z)(y - x)\|_\infty \quad (19)$$

$$\stackrel{(c)}{\leq} \sqrt{n} n \|\nabla \text{Softmax}^{(n)}(z)\|_\infty \|y - x\|_\infty \quad (20)$$

$$\stackrel{(d)}{\leq} n^{3/2} \|y - x\|_\infty, \quad (21)$$

where (a) and (b) follow from Lemma 2 in the Handout, (c) is by Lemma 3 in the Handout, and (d) is by subproblem (c). We conclude that

$$|\text{Softmax}^{(n)}|_{Lip} := \sup_{\substack{x, y \in [-1, 1]^d \\ x \neq y}} \frac{\|\text{Softmax}^{(n)}(x) - \text{Softmax}^{(n)}(y)\|_\infty}{\|x - y\|_\infty} \leq n^{3/2}. \quad (22)$$

(e) Let  $x, y \in [-1, 1]^d$ .

$$\|\phi(x) - \phi(y)\|_\infty = \|\text{Softmax}^{(n)}(A_2 \rho(A_1 x)) - \text{Softmax}^{(n)}(A_2 \rho(A_1 y))\|_\infty \quad (23)$$

$$\stackrel{(a)}{\leq} n^{3/2} \|A_2 \rho(A_1 x) - A_2 \rho(A_1 y)\|_\infty \quad (24)$$

$$\stackrel{(b)}{\leq} n^{3/2} m \|A_2\|_\infty \|\rho(A_1 x) - \rho(A_1 y)\|_\infty \quad (25)$$

$$\stackrel{(c)}{\leq} n^{3/2} m \|A_2\|_\infty \|A_1 x - A_1 y\|_\infty \quad (26)$$

$$\stackrel{(d)}{\leq} n^{3/2} m d \|A_2\|_\infty \|A_1\|_\infty \|x - y\|_\infty, \quad (27)$$

where (a) is by subproblem (d), (b) follows from Lemma 3 in the Handout, (c) is a consequence of the ReLU function being 1-Lipschitz, and (d) follows from Lemma 3 in the Handout. In summary, we have established that

$$|\phi|_{Lip} \leq n^{3/2} m d \|A_2\|_\infty \|A_1\|_\infty. \quad (28)$$

## Problem 4

(a) We use Definition 6 in the Handout to compute  $x_1, x_2, x_3$ , and  $x_4$  as follows:

$$x_1 = (A_1 * K_1)_{3,2} = \sum_{\substack{p \in \{1, \dots, 1\} \\ q \in \{1, \dots, 3\}}} [A_1]_{3+p-1, 2+q-1} [K_1]_{p,q} = \sum_{q \in \{1, \dots, 3\}} [A_1]_{3, 1+q} [K_1]_{1,q} \quad (29)$$

$$= [A_1]_{3,2} [K_1]_{1,1} + [A_1]_{3,3} [K_1]_{1,2} + [A_1]_{3,4} [K_1]_{1,3} = 1. \quad (30)$$

$$x_2 = (A_1 * K_1)_{4,2} = \sum_{\substack{p \in \{1, \dots, 1\} \\ q \in \{1, \dots, 3\}}} [A_1]_{4+p-1, 2+q-1} [K_1]_{p,q} = \sum_{q \in \{1, \dots, 3\}} [A_1]_{4, 1+q} [K_1]_{1,q} \quad (31)$$

$$= [A_1]_{4,2} [K_1]_{1,1} + [A_1]_{4,3} [K_1]_{1,2} + [A_1]_{4,4} [K_1]_{1,3} = 1. \quad (32)$$

$$x_3 = (A_1 * K_2)_{2,4} = \sum_{\substack{p \in \{1, \dots, 3\} \\ q \in \{1, \dots, 1\}}} [A_1]_{2+p-1, 4+q-1} [K_2]_{p,q} = \sum_{p \in \{1, \dots, 3\}} [A_1]_{1+p, 4} [K_2]_{p,1} \quad (33)$$

$$= [A_1]_{2,4} [K_2]_{1,1} + [A_1]_{3,4} [K_2]_{2,1} + [A_1]_{4,4} [K_2]_{3,1} = 2. \quad (34)$$

$$x_4 = (A_1 * K_2)_{3,4} = \sum_{\substack{p \in \{1, \dots, 3\} \\ q \in \{1, \dots, 1\}}} [A_1]_{3+p-1, 4+q-1} [K_2]_{p,q} = \sum_{p \in \{1, \dots, 3\}} [A_1]_{2+p, 4} [K_2]_{p,1} \quad (35)$$

$$= [A_1]_{3,4} [K_2]_{1,1} + [A_1]_{4,4} [K_2]_{2,1} + [A_1]_{5,4} [K_2]_{3,1} = 3. \quad (36)$$

We deduce that

$$\|A_1 * K_1\|_\infty = 1, \quad \|A_1 * K_2\|_\infty = 3. \quad (37)$$

It follows that

$$\phi(A_1) = (\|A_1 * K_1\|_\infty, \|A_1 * K_2\|_\infty) = (1, 3). \quad (38)$$

(b) As  $\phi(A_1) = \phi(A_2)$ ,  $X$  is not in  $\phi$ -general position.  $\{X^+, X^-\}$  is  $\phi$ -separable. Indeed, consider  $w = (-1, 1)$ . Then,

$$\langle \phi(A_1), w \rangle = \langle \phi(A_2), w \rangle = 2, \quad (39)$$

and

$$\langle \phi(A_3), w \rangle = \langle \phi(A_4), w \rangle = -2. \quad (40)$$

(c)  $\|K_1\|_1 = \|K_2\|_1 = \frac{1}{11} \sum_{i=1}^{11} 1 = 1$ . By Lemma 4 in the Handout, for all  $A \in [0, 1]^{28 \times 28}$ ,  $\|A * K_1\|_\infty \leq \|A\|_\infty \|K_1\|_1 \leq \|A\|_\infty$  and  $\|A * K_2\|_\infty \leq \|A\|_\infty \|K_2\|_1 \leq \|A\|_\infty$ . As  $A \in [0, 1]^{28 \times 28}$ , it follows that  $\phi(A) \in [0, 1]^2$ , for all  $A \in X$ .

(d) In subproblem (a), we have seen that convolution with a matrix of ones lined up horizontally tends to produce maximal values if the image contains a horizontal line, and conversely, convolution with a matrix of ones lined up vertically tends to produce maximal values if the image contains a vertical line.

The red points represent the value of  $\phi(A)$  for  $A \in X^-$ .  $X^-$  contains images of handwritten ones, which are approximately vertical lines. Hence, for  $A \in X^-$ ,



$A * K_1$  yields small values (close to 0) while  $A * K_2$  yields higher values (close to 1), and therefore  $\|A * K_1\|_\infty \simeq 0$  while  $\|A * K_2\|_\infty \simeq 1$ . This explains why  $\phi(A) \simeq (0, 1)$  for  $A \in X^-$ .

The blue points represent the value of  $\phi(A)$  for  $A \in X^+$ .  $X^+$  contains images of handwritten zeros, which are approximately a combination of horizontal and vertical lines. Hence, for  $A \in X^+$ , both  $A * K_1$  and  $A * K_2$  yield high values (close to 1), and therefore  $\|A * K_1\|_\infty \simeq \|A * K_2\|_\infty \simeq 1$ . This explains why  $\phi(A) \simeq (1, 1)$  for  $A \in X^+$ .

- (e) The distribution of blue and red points clearly overlap, so it is not possible to separate  $\{\phi(X^+), \phi(X^-)\}$  with a line. Therefore,  $X$  is not  $\phi$ -separable. We can design a map  $\phi'$  which could lead to better separation properties as follows.  $\phi$  is designed to “spot” the horizontal and vertical lines on images, through the convolution with  $K_1$  and  $K_2$ . However, handwritten zeros and ones are not composed of horizontal and vertical lines: ones are not perfectly vertical, but rather deviate a little bit to the right, and zeros have many components, such as curved lines and edges. A strategy to get a better separation would be to add more coordinates to the output of  $\phi$  to detect more diverse features. For example, one can choose  $\phi' : \mathbb{R}^{28 \times 28} \rightarrow \mathbb{R}^{12}$ , according to

$$\phi'(A) = (A * K_i)_{i \in \{1, \dots, 12\}}, \quad \forall A \in \mathbb{R}^{28 \times 28}, \quad (41)$$

where the  $K_i$  are matrices. The first two coordinates would remain the same, but the next ones would use convolutions with other matrices, with the purpose of detecting inclined lines, curved lines and edges. Here is an example of the  $K_i$ , for  $i = 3, \dots, 12$ :

$$K_3 = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}, \quad K_4 = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \quad (42)$$

are designed to detect inclined lines,

$$K_5 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad K_6 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad K_7 = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 \end{pmatrix}, \quad K_8 = \begin{pmatrix} 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix}$$

are designed to detect curved lines, and

$$K_9 = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}, \quad K_{10} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad K_{11} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad K_{12} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \quad (43)$$

are designed to detect edges.