

Solutions to the Exam on Neural Network Theory February 6, 2020

Problem 1

(a) We have to establish that $K(\boldsymbol{x}, \boldsymbol{y}) = K(\boldsymbol{y}, \boldsymbol{x})$, which follows from

$$\begin{split} K(\boldsymbol{x},\boldsymbol{y}) &= \alpha x_1^2 y_1^2 + \beta x_2^2 y_2^2 \\ &= \alpha y_1^2 x_1^2 + \beta y_2^2 x_2^2 \\ &= K(\boldsymbol{y},\boldsymbol{x}), \quad \text{for all } \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^2 \end{split}$$

(b) Let $\varepsilon > 0$ and $x, y \in \mathbb{R}^2$. We have to show that there exists a $\delta = \delta(\varepsilon, x, y) > 0$ such that if $u, v \in \mathbb{R}^2$ satisfy

$$\left\| \begin{pmatrix} \boldsymbol{x} \\ \boldsymbol{y} \end{pmatrix} - \begin{pmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{pmatrix} \right\|_2 < \delta, \tag{1}$$

then $|K(\boldsymbol{x}, \boldsymbol{y}) - K(\boldsymbol{u}, \boldsymbol{v})| < \varepsilon$. Since

$$\begin{aligned} |K(\boldsymbol{x}, \boldsymbol{y}) - K(\boldsymbol{u}, \boldsymbol{v})| &= \left| \alpha x_1^2 y_1^2 + \beta x_2^2 y_2^2 - \alpha u_1^2 v_1^2 - \beta u_2^2 v_2^2 \right| \\ &\leq \alpha \left| x_1^2 y_1^2 - u_1^2 v_1^2 \right| + \beta \left| x_2^2 y_2^2 - u_2^2 v_2^2 \right|, \end{aligned}$$

it is sufficient to choose $\delta > 0$ such that

$$\left|x_1^2y_1^2-u_1^2v_1^2\right|<\frac{\varepsilon}{2\alpha}$$

and

$$\left|x_2^2 y_2^2 - u_2^2 v_2^2\right| < \frac{\varepsilon}{2\beta}.$$

Suppose that (1) holds. This implies that

 $|x_1 - u_1| < \delta$, $|y_1 - v_1| < \delta$, $|x_2 - u_2| < \delta$, and $|y_2 - v_2| < \delta$.

It now follows that

$$\begin{aligned} \left| x_1^2 y_1^2 - u_1^2 v_1^2 \right| &= \left| x_1^2 y_1^2 - x_1^2 v_1^2 + x_1^2 v_1^2 - u_1^2 v_1^2 \right| \\ &\leq x_1^2 \left| y_1^2 - v_1^2 \right| + v_1^2 \left| x_1^2 - u_1^2 \right| \\ &= x_1^2 \left| (y_1 + v_1) (y_1 - v_1) \right| + v_1^2 \left| (x_1 + u_1) (x_1 - u_1) \right| \\ &\leq x_1^2 (|y_1| + |v_1|) \delta + v_1^2 (|x_1| + |u_1|) \delta \\ &\leq \left(x_1^2 (2|y_1| + \delta) + (|y_1| + \delta)^2 (2|x_1| + \delta) \right) \delta \\ &< \frac{\varepsilon}{2\alpha} \end{aligned}$$

provided that $\delta < C_1$ with

$$C_1 = \min\left\{1, \frac{\varepsilon}{2\alpha(x_1^2(2|y_1|+1) + (|y_1|+1)^2(2|x_1|+1))}\right\}.$$

Swapping the roles of x_1 and x_2 , y_1 and y_2 , u_1 and u_2 , v_1 and v_2 , and α and β and using the same line of arguments yields

$$\left|x_2^2 y_2^2 - u_2^2 v_2^2\right| < \frac{\varepsilon}{2\beta}$$

provided that $\delta < C_2$ with

$$C_2 = \min\left\{1, \frac{\varepsilon}{2\beta(x_2^2(2|y_2|+1) + (|y_2|+1)^2(2|x_2|+1))}\right\}.$$

We conclude that if $\delta < \min\{C_1, C_2\}$, then

$$\left\| \begin{pmatrix} \boldsymbol{x} \\ \boldsymbol{y} \end{pmatrix} - \begin{pmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{pmatrix} \right\|_2 < \delta$$

implies $|K(\boldsymbol{x}, \boldsymbol{y}) - K(\boldsymbol{u}, \boldsymbol{v})| < \varepsilon$, which establishes continuity of $K(\boldsymbol{x}, \boldsymbol{y})$.

(c) We have to show that for every $k \in \mathbb{N}$ and $x_1, \ldots, x_k \in \mathbb{R}^2$, the $k \times k$ Gramian matrix

$$\boldsymbol{K}(\boldsymbol{x}_1,\ldots,\boldsymbol{x}_k) = \begin{pmatrix} K(\boldsymbol{x}_1,\boldsymbol{x}_1) & K(\boldsymbol{x}_1,\boldsymbol{x}_2) & \ldots & K(\boldsymbol{x}_1,\boldsymbol{x}_k) \\ K(\boldsymbol{x}_2,\boldsymbol{x}_1) & K(\boldsymbol{x}_2,\boldsymbol{x}_2) & \ldots & K(\boldsymbol{x}_2,\boldsymbol{x}_k) \\ \vdots & \vdots & \ddots & \vdots \\ K(\boldsymbol{x}_k,\boldsymbol{x}_1) & K(\boldsymbol{x}_k,\boldsymbol{x}_2) & \ldots & K(\boldsymbol{x}_k,\boldsymbol{x}_k) \end{pmatrix}$$

is positive semidefinite. To this end, consider the mapping $\Phi\colon\mathbb{R}^2\to\mathbb{R}^2$ defined as

$$\Phi(\boldsymbol{x}) = \begin{pmatrix} \sqrt{\alpha}x_1^2\\ \sqrt{\beta}x_2^2 \end{pmatrix}$$

and note that $K(\boldsymbol{x}, \boldsymbol{y}) = \Phi^{\mathsf{T}}(\boldsymbol{x})\Phi(\boldsymbol{y})$. It follows that, for every $k \in \mathbb{N}$ and $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_k \in \mathbb{R}^2$, we have

$$\boldsymbol{c}^{\mathsf{T}}\boldsymbol{K}(\boldsymbol{x}_{1},\ldots,\boldsymbol{x}_{k})\boldsymbol{c} = \sum_{i,j=1}^{k} c_{i}c_{j}K(\boldsymbol{x}_{i},\boldsymbol{x}_{j})$$
$$= \sum_{i,j=1}^{k} c_{i}c_{j}\Phi^{\mathsf{T}}(\boldsymbol{x}_{i})\Phi(\boldsymbol{x}_{j})$$
$$= \left(\sum_{i=1}^{k} c_{i}\Phi(\boldsymbol{x}_{i})\right)^{\mathsf{T}}\left(\sum_{j=1}^{k} c_{j}\Phi(\boldsymbol{x}_{j})\right)$$
$$\geq 0, \quad \text{for all } \boldsymbol{c} = (c_{1}\ldots c_{k})^{\mathsf{T}} \in \mathbb{R}^{k},$$

which proves that K(x, y) is positive semidefinite.

(d) For every $y \in \mathbb{R}^2$, we define the function $K_y \colon \mathbb{R}^2 \to \mathbb{R}$ according to $K_y(x) = K(x, y)$ and set

$$\mathcal{H}_0 = \operatorname{span}\{K_{\boldsymbol{y}} : \boldsymbol{y} \in \mathbb{R}^2\}.$$

Let $e_1 = (1 \ 0)^{\mathsf{T}}$ and $e_2 = (0 \ 1)^{\mathsf{T}}$ with corresponding

$$K_{\boldsymbol{e}_1}(\boldsymbol{x}) = K(\boldsymbol{e}_1, \boldsymbol{x}) = \alpha x_1^2$$

and

$$K_{\boldsymbol{e}_2}(\boldsymbol{x}) = K(\boldsymbol{e}_2, \boldsymbol{x}) = \beta x_2^2,$$

respectively. Since

$$egin{aligned} K_{oldsymbol{a}}(oldsymbol{x}) &= lpha x_1^2 a_1^2 + eta x_2^2 a_2^2 \ &= a_1^2 K_{oldsymbol{e}_1}(oldsymbol{x}) + a_2^2 K_{oldsymbol{e}_2}(oldsymbol{x}), & ext{ for all } oldsymbol{a} &= egin{aligned} a_1 \ a_2 \end{pmatrix} \in \mathbb{R}^2, \end{aligned}$$

it follows that $\mathcal{H}_0 = \operatorname{span}\{K_{e_1}, K_{e_2}\}$. We conclude that the reproducing kernel Hilbert space \mathcal{H}_K corresponding to the kernel $K(\boldsymbol{x}, \boldsymbol{y})$ is given by

$$\mathcal{H}_K = \overline{\mathcal{H}_0} = \overline{\operatorname{span}\{K_{\boldsymbol{e}_1}, K_{\boldsymbol{e}_2}\}} = \operatorname{span}\{K_{\boldsymbol{e}_1}, K_{\boldsymbol{e}_2}\}.$$

(e) Consider K_{e_1} and K_{e_2} from subproblem (d). Using the reproducing property of \mathcal{H}_K , we have

$$\langle K_{\boldsymbol{e}_1}, K_{\boldsymbol{e}_1} \rangle_{\mathcal{H}_K} = K(\boldsymbol{e}_1, \boldsymbol{e}_1) = \alpha, \langle K_{\boldsymbol{e}_2}, K_{\boldsymbol{e}_2} \rangle_{\mathcal{H}_K} = K(\boldsymbol{e}_2, \boldsymbol{e}_2) = \beta,$$

and

$$\langle K_{\boldsymbol{e}_1}, K_{\boldsymbol{e}_2} \rangle_{\mathcal{H}_K} = K(\boldsymbol{e}_1, \boldsymbol{e}_2) = 0.$$

This implies that $\{K_{e_1}/\sqrt{\alpha}, K_{e_2}/\sqrt{\beta}\}$ is an orthonormal basis for the reproducing kernel Hilbert space $\mathcal{H}_K = \operatorname{span}\{K_{e_1}, K_{e_2}\}$ corresponding to the kernel $K(\boldsymbol{x}, \boldsymbol{y})$.

Problem 2

(a) The following figure depicts $\{x_1, x_2\}$ (red points) and $\{x_3, x_4\}$ (blue points) together with a separating straight line (thick black line):



(b) The Lagrange function L(w, b, c) id given by

$$L(\boldsymbol{w}, b, \boldsymbol{c}) = \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + \sum_{i=1}^{4} c_{i}(1 - y_{i}(\langle \boldsymbol{w}, \boldsymbol{x}_{i} \rangle_{2} - b)).$$
(2)

(c) Setting $\nabla_{\boldsymbol{w}} L(\boldsymbol{w}, b, \boldsymbol{c}) = 0$ and $\nabla_{b} L(\boldsymbol{w}, b, \boldsymbol{c}) = 0$ yields

$$\boldsymbol{w} = \sum_{i=1}^{4} c_i y_i \boldsymbol{x}_i \tag{3}$$

and

$$\sum_{i=1}^{4} c_i y_i = 0,$$
(4)

respectively. Using (3) and (4) in (2) results in

$$g(\boldsymbol{c}) = \min_{\boldsymbol{w} \in \mathbb{R}^2, b \in \mathbb{R}} L(\boldsymbol{w}, b, \boldsymbol{c})$$
$$= \sum_{i=1}^4 c_i - \frac{1}{2} \sum_{i,j=1}^4 c_i y_i c_j y_j \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle_2$$
$$= \boldsymbol{a}^{\mathsf{T}} \boldsymbol{c} - \frac{1}{2} \boldsymbol{c}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{c}$$

with $\boldsymbol{a} = (1 \ 1 \ 1 \ 1)^{\mathsf{T}}$ and

$$\boldsymbol{A} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

(d) It follows from subproblem (c) that the Lagrange dual function g(c) can be written as

$$g(\mathbf{c}) = c_1 + c_2 + c_3 + c_4 - \frac{1}{2}(c_2^2 + c_3^2 + 2c_2c_3 + c_4^2).$$

Using the constraint

$$\sum_{i=1}^{4} c_i y_i = c_1 + c_2 - c_3 - c_4 = 0,$$
(5)

we can write

$$g(\mathbf{c}) = 2(c_3 + c_4) - \frac{1}{2}(c_3^2 + c_4^2) - \frac{1}{2}(c_2^2 + 2c_2c_3).$$

It follows that

$$g(\mathbf{c}) \le 2(c_3 + c_4) - \frac{1}{2}(c_3^2 + c_4^2)$$

$$\le 4$$

with equalities in both steps for $\tilde{c}_2 = 0$ and $\tilde{c}_3 = \tilde{c}_4 = 2$. Using (5) again, we get $\tilde{c}_1 = \tilde{c}_3 + \tilde{c}_4 - \tilde{c}_2 = 4$. The solution of the Lagrange dual problem is therefore given by $\tilde{c} = (4 \ 0 \ 2 \ 2)^T$. A vector x_i is a support vector if and only if the corresponding \tilde{c}_i is strictly positive. The support vectors are therefore given by $\{x_1, x_3, x_4\}$.

(e) Using (3) and the solution $\tilde{c} = (4 \ 0 \ 2 \ 2)^{\mathsf{T}}$ from subproblem (d), we obtain the following solution \tilde{w} of the optimization problem in subproblem (b):

$$\tilde{\boldsymbol{w}} = 4\boldsymbol{x}_1 - 2\boldsymbol{x}_3 - 2\boldsymbol{x}_4 = -2\begin{pmatrix}1\\1\end{pmatrix}.$$

The solution \tilde{b} can be obtained from

$$b = \langle \tilde{\boldsymbol{w}}, \boldsymbol{x}_1 \rangle_2 - y_1 \\ = -1,$$

which follows from the fact that for the support vector \boldsymbol{x}_1 we have $\tilde{c}_1 > 0$ and, therefore, the corresponding inequality constraint $y_1(\langle \tilde{\boldsymbol{w}}, \boldsymbol{x}_1 \rangle_2 - \tilde{b}) \geq 1$ must be satisfied with equality. The hard margin binary classifier $g_{\text{hm}}(\boldsymbol{x})$ therefore has the following form:

$$g_{\mathsf{hm}}(\boldsymbol{x}) = (\langle \tilde{\boldsymbol{w}}, \boldsymbol{x} \rangle_2 - b) \ = -2 \begin{pmatrix} 1 & 1 \end{pmatrix} \boldsymbol{x} + 1.$$

Problem 3

(a) The following figure depicts f(x).



(b) Denote the ReLU function as $\rho(x) = \max\{0, x\}$. The function f(x) can be realized as a linear combination of shifted ReLU functions according to

$$f(x) = \rho(4x) - \rho(4x - 1) - \rho(4x - 3) + \rho(4x - 4).$$

This function can be realized through a depth-2 ReLU network according to

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

with

$$W_1(x) = \begin{pmatrix} 4\\4\\4\\4 \end{pmatrix} x + \begin{pmatrix} 0\\-1\\-3\\-4 \end{pmatrix}, \quad W_2(x) = \begin{pmatrix} 1 & -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1\\x_2\\x_3\\x_4 \end{pmatrix}.$$

The network Φ has depth 2, width 4, and connectivity 11.

(c) One way to realize h(x) through a ReLU-network is according to

$$h(x) = f(4x) + f(4x - 3) = \rho(16x) - \rho(16x - 1) - \rho(16x - 3) + \rho(16x - 4) + \rho(16x - 12) - \rho(16x - 13) - \rho(16x - 15) + \rho(16x - 16).$$

This corresponds to the depth-2 ReLU-network

$$\Phi_1(x) = W_2(\rho(W_1(x)))$$

with

$$W_{1}(x) = \begin{pmatrix} 16\\16\\16\\16\\16\\16\\16\\16\\16\\16 \end{pmatrix} x + \begin{pmatrix} 0\\-1\\-3\\-4\\-12\\-13\\-15\\-16 \end{pmatrix}, \quad W_{2}(x) = \begin{pmatrix} 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_{1}\\x_{2}\\x_{3}\\x_{4}\\x_{5}\\x_{6}\\x_{7}\\x_{8} \end{pmatrix}.$$

The network Φ_1 has depth 2, width 8, and connectivity 23.

(d) An alternative way to realize h(x) is to note that h(x) = f(f(x)). The function f(x) is realized through the network $\Phi(x) = W_2(\rho(W_1(x)))$ in subproblem (b), therefore h(x) = f(f(x)) can be realized through the network

$$\Phi_2(x) = \Phi(\Phi(x)) = W_2(\rho(W_1(W_2(\rho(W_1(x)))))) = W'_3(\rho(W'_2(\rho(W'_1(x)))))$$

with

The network Φ_2 has depth 3, width 4, and connectivity 30.

Problem 4

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

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

(b) Let $\{x_1, \ldots, x_N\}$ be an ϵ -packing of C. The ℓ_2 -balls centered at the x_i and of radius $\epsilon/2$ are referred to as packing balls. As the distance between each pair in the packing is at least ϵ , the packing balls are mutually disjoint. Therefore, the maximum volume covered by the packing balls is $M(\epsilon; C, \|\cdot\|_2) c(\epsilon/2)^d$, where c is a constant. Since each point in the packing must be in C, all the packing balls are, indeed, contained in an ℓ_2 -ball of radius $1 + \epsilon/2$. Hence, we get, for all $\epsilon \leq 2$,

$$M(\epsilon; \mathcal{C}, \|\cdot\|_2) \ c \ (\epsilon/2)^d \le c \ (1+\epsilon/2)^d \Longrightarrow M(\epsilon; \mathcal{C}, \|\cdot\|_2) \le \left(\frac{2}{\epsilon}+1\right)^d \le \left(\frac{4}{\epsilon}\right)^d.$$