

# Kernel Methods AMMI 2023

## Final Exam

Juliette Marrie and Jean-Philippe Vert

July 2023

### 1 Dual of ridge regression (8.5 points)

Given  $x_1, \dots, x_n \in \mathbb{R}^p$  and  $y_1, \dots, y_n \in \mathbb{R}$ , we consider the ridge regression problem

$$\min_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|^2$$

**Question 1:** ..... 2 points

- (a) (.5 points) Rewrite this problem in matrix notation, with the matrix  $X \in \mathbb{R}^{n \times p}$  and vector  $y \in \mathbb{R}^n$  defined by  $X = [x_1, \dots, x_n]^\top$  and  $y = (y_1, \dots, y_n)^\top$ .

**Solution:**

$$\min_{w \in \mathbb{R}^p} \frac{1}{n} (Xw - y)^\top (Xw - y) + \lambda w^\top w$$

*Note: Those who wrote  $wX$  instead of  $Xw$  did not get the points.*

- (b) (1.5 points) Find  $w^* \in \mathbb{R}^p$  that solves the problem.

**Solution:**

$$w = (X^\top X + \lambda n I)^{-1} X^\top y \quad (\text{proof in slides})$$

**Question 2:** ..... 6 points

Let  $u = Xw$ . We rewrite the ridge regression problem as

$$\begin{aligned} \min_{u \in \mathbb{R}^n, w \in \mathbb{R}^p} \quad & \frac{1}{n} \|u - y\|^2 + \lambda \|w\|^2 \\ \text{s.t.} \quad & u = Xw \end{aligned}$$

- (a) (.5 points) Write the Lagrangian of this primal problem.

**Solution:**

$$\mathcal{L}(w, u, \alpha) = \frac{1}{n} \|u - y\|^2 + \lambda \|w\|^2 + \alpha^\top (u - Xw), \quad \alpha \in \mathbb{R}^n$$

- (b) (3 points) Show that the dual problem can be written as

$$\max_{\alpha \in \mathbb{R}^n} -\frac{1}{4\lambda} [\alpha^\top (XX^\top + \lambda n I) \alpha - 4\lambda \alpha^\top y]$$

**Solution:** The dual problem is

$$\max_{\alpha \in \mathbb{R}^n} q(\alpha) \quad \text{with} \quad q(\alpha) = \inf_{w, u} \mathcal{L}(w, u, \alpha)$$

$$\nabla_w \mathcal{L} = 2\lambda w - X^T \alpha = 0 \implies w = \frac{X^T \alpha}{2\lambda}$$

$$\nabla_u \mathcal{L} = \frac{2}{n}(u - y) + \alpha = 0 \implies u = y - \frac{n\alpha}{2}$$

Therefore,

$$\begin{aligned} q(\alpha) &= \frac{1}{n} \left\| y - \frac{n\alpha}{2} - y \right\|^2 + \lambda \left\| \frac{X^T \alpha}{2\lambda} \right\|^2 + \alpha^T \left( y - \frac{n\alpha}{2} - \frac{XX^T \alpha}{2\lambda} \right) \\ &= \frac{n}{4} \alpha^T \alpha + \frac{1}{4\lambda} \alpha^T XX^T \alpha + \alpha^T y - \frac{n}{2} \alpha^T \alpha - \frac{1}{2\lambda} \alpha^T XX^T \alpha \\ &= -\frac{n}{4} \alpha^T \alpha - \frac{1}{4\lambda} \alpha^T XX^T \alpha + \alpha^T y \\ &= -\frac{1}{4\lambda} [\alpha^T (XX^T + n\lambda I) \alpha - 4\lambda \alpha^T y] \end{aligned}$$

- (c) (1 point) Find  $\alpha^*$  that solves this problem.

**Solution:**

$$\nabla_{\alpha} q = 2(XX^T + n\lambda I)\alpha - 4\lambda y = 0 \implies \alpha = 2\lambda(XX^T + n\lambda I)^{-1}y$$

- (d) (.5 points) Deduce the value of  $w^*$  that solves the primal problem.

**Solution:**

$$w = \frac{X^T \alpha}{2\lambda} = X^T (XX^T + \lambda n I)^{-1} y$$

- (e) (1.5 points) Show that solutions  $w^*$  found in Question 1(b) and in Question 2(d) are equal.

**Solution:** We want to show that

$$X^T (XX^T + \lambda n I)^{-1} y = (X^T X + \lambda n I)^{-1} X^T y$$

We have

$$X^T XX^T + \lambda n I_p X^T = X^T XX^T + \lambda n X^T I_n$$

Factorizing by  $X^T$ ,

$$(X^T X + \lambda n I_p) X^T = X^T (XX^T + \lambda n I_n)$$

Multiplying by  $(X^T X + \lambda n I_p)^{-1}$  on the left and  $(XX^T + \lambda n I_n)^{-1}$  on the right, on both sides of the equality, we get

$$X^T (XX^T + \lambda n I_n)^{-1} = (X^T X + \lambda n I_p)^{-1} X^T$$

## 2 Some p.d. kernels (3 points)

Show that the following kernels are p.d. (hint: write them as inner products):

a) (1.5 point)  $K : \begin{cases} \mathbb{R}^2 & \rightarrow \\ (x, y) & \mapsto \end{cases} K(x, y) = 2^{x+y} + 3^{x+y}$

**Solution:**

$$K(x, y) = \phi(x)^T \phi(y) \quad \text{with} \quad \phi(x) = (2^x, 3^x)^T$$

b) (1.5 point)  $K : \begin{cases} \mathbb{R}^2 & \rightarrow \\ (x, y) & \mapsto \end{cases} K(x, y) = \cos(x - y)$

**Solution:** We have  $\cos(x - y) = \cos(x)\cos(y) + \sin(x)\sin(y)$ , therefore

$$K(x, y) = \phi(x)^T \phi(y) \quad \text{with} \quad \phi(x) = (\cos(x), \sin(x))^T$$

## 3 Support Vector Regression (SVR) (5.5 pointw)

Given  $x_1, \dots, x_n \in \mathbb{R}^p$  and  $y_1, \dots, y_n \in \mathbb{R}$ , linear SVR finds a linear model  $f(x) = w^\top x + b$  with the maximum number of points such as the prediction  $f(x_i)$  is within  $\pm\epsilon$  of the output  $y_i$ . Samples with prediction error at least  $\epsilon$  penalize the objective by  $|f(x_i) - y_i| - \epsilon$ .

The SVR problem can be formulated as the following optimization problem:

$$\begin{aligned} \min_{w \in \mathbb{R}^p, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} w^\top w + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i - w^\top x_i - b \leq \epsilon + \xi_i \\ & y_i - w^\top x_i - b \geq -\epsilon - \xi_i \\ & \xi_i \geq 0, i = 1, \dots, n \end{aligned}$$

Rewrite this problem as a quadratic program, i.e., find  $z, P, q, G$  and  $h$  such that the problem can be formulated as

$$\begin{aligned} \min_z \quad & z^\top P z + q^\top z \\ \text{s.t.} \quad & G z \leq h \end{aligned}$$

In other words,

- (.5 points) define the vector  $z$  and give its dimension
- (1 point) give the dimensions of  $P, q, G$  and  $h$
- (1 + .5 + 1.5 + 1 = 4 points) give their value: you can write them as block vectors or block matrices, indicating the dimension of each block

**Solution:** We rewrite the SVR problem in the QP form:

$$\begin{aligned} \min_{w \in \mathbb{R}^p, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} w^\top w + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & -w^\top x_i - b - \xi_i \leq \epsilon - y_i \\ & w^\top x_i + b - \xi_i \leq \epsilon + y_i \\ & -\xi_i \leq 0, \quad i = 1, \dots, n \end{aligned}$$

We can define  $z = (w_1, \dots, w_p, b, \xi_1, \dots, \xi_n)^\top \in \mathbb{R}^{p+1+n}$ .

Let  $m = p + 1 + n$ . The shapes of  $P, q, G$  and  $h$  are respectively  $(m, m)$ ,  $(m, )$ ,  $(3n, m)$  and  $(3n, )$ .

- $P = \text{diag}(1, \dots, 1, 0, 0, \dots, 0)$  with  $p$  ones and  $1 + n$  zeros.
- $q = (0, \dots, 0, 0, C, \dots, C)^\top$  with  $p + 1$  zeros and  $n$   $C$ .

$$\bullet \quad G = \left( \begin{array}{c|c|c} \mathbf{-X} & \begin{matrix} -1 \\ \vdots \\ -1 \end{matrix} & \begin{matrix} -1 & & \\ & \ddots & \\ & & -1 \end{matrix} \\ \hline \mathbf{X} & \begin{matrix} 1 \\ \vdots \\ 1 \end{matrix} & \begin{matrix} -1 & & \\ & \ddots & \\ & & -1 \end{matrix} \\ \hline \mathbf{0} & \begin{matrix} 0 \\ \vdots \\ 0 \end{matrix} & \begin{matrix} -1 & & \\ & \ddots & \\ & & -1 \end{matrix} \end{array} \right)$$

- $h = (\epsilon - y_1, \dots, \epsilon - y_n, \epsilon + y_1, \dots, \epsilon + y_n, 0, \dots, 0)^\top$ .

## 4 Maximum Mean Discrepancy (3 points)

Let  $\mathcal{X}$  be a set and  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$  a RKHS with reproducing kernel  $K$ .

Given  $n$  points  $X = (x_1, \dots, x_n) \in \mathcal{X}^n$  we define

$$\hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n K_{x_i},$$

where for any  $x \in \mathcal{X}$ ,  $K_x \in \mathcal{H}$  denotes the function  $t \in \mathcal{X} \mapsto K(x, t) \in \mathbb{R}$ .

**Question 1:** ..... 1.5 points

We consider two sets of points  $X = (x_1, \dots, x_n) \in \mathcal{X}^n$  and  $Y = (y_1, \dots, y_m) \in \mathcal{X}^m$ . Express  $\|\hat{\mu}_X - \hat{\mu}_Y\|_{\mathcal{H}}^2$  as a function of  $K$ .

**Solution:**

$$\begin{aligned}
\|\hat{\mu}_X - \hat{\mu}_Y\|_{\mathcal{H}}^2 &= \langle \hat{\mu}_X, \hat{\mu}_X \rangle_{\mathcal{H}} - 2\langle \hat{\mu}_X, \hat{\mu}_Y \rangle_{\mathcal{H}} + \langle \hat{\mu}_Y, \hat{\mu}_Y \rangle_{\mathcal{H}} \\
&= \left\langle \frac{1}{n} \sum_{i=1}^n K_{x_i}, \frac{1}{n} \sum_{j=1}^n K_{x_j} \right\rangle_{\mathcal{H}} - 2 \left\langle \frac{1}{n} \sum_{i=1}^n K_{x_i}, \frac{1}{m} \sum_{j=1}^m K_{y_j} \right\rangle_{\mathcal{H}} + \left\langle \frac{1}{m} \sum_{i=1}^m K_{y_i}, \frac{1}{m} \sum_{j=1}^m K_{y_j} \right\rangle_{\mathcal{H}} \\
&= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \langle K_{x_i}, K_{x_j} \rangle_{\mathcal{H}} - 2 \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \langle K_{x_i}, K_{y_j} \rangle_{\mathcal{H}} + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \langle K_{y_i}, K_{y_j} \rangle_{\mathcal{H}} \\
&= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(x_i, x_j) - 2 \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m K(x_i, y_j) + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m K(y_i, y_j)
\end{aligned}$$

**Question 2:** ..... 1.5 points

For any  $f \in \mathcal{H}$ , express  $\langle f, \hat{\mu}_X \rangle_{\mathcal{H}}$  as a function of  $f(x_1), \dots, f(x_n)$ .

**Solution:**

$$\langle f, \hat{\mu}_X \rangle_{\mathcal{H}} = \left\langle f, \frac{1}{n} \sum_{i=1}^n K_{x_i} \right\rangle_{\mathcal{H}} = \frac{1}{n} \sum_{i=1}^n \langle f, K_{x_i} \rangle_{\mathcal{H}} = \sum_{i=1}^n f(x_i) \quad \text{since} \quad \langle f, K_x \rangle_{\mathcal{H}} = f(x)$$

*Note:  $f$  is not the solution of an optimization problem depending on  $x_1, \dots, x_n$ . There is no supposed link between  $f$  and  $x_1, \dots, x_n$ , therefore  $f$  does not necessarily lie the span of the  $K_{x_i}$ .*