

# Machine Learning 2024 - Sheet 3.2

## Block III: SVM and Kernel Methods

Isabel Valera

### Exercise 1: Kernel feature representation



Given the kernel  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + c)^2$ . Write down, step-by-step, a feature representation  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$  where  $d < p$  such that  $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = k(\mathbf{x}, \mathbf{z})$ .

### Exercise 2: Eigenvalues



**Definition Positive Definite Matrix** A complex  $m \times m$  matrix  $K$  satisfying

$$\sum_{i,j} c_i \bar{c}_j K_{ij} \geq 0 \quad (1)$$

for all  $c_i \in \mathbb{C}$  is called positive definite. The bar in  $\bar{c}_j$  denotes complex conjugation; for real numbers, it has no effect. Similarly, a real symmetric  $m \times m$  matrix  $K$  satisfying (1) for all  $c_i \in \mathbb{R}$  is called positive definite.

Prove that a symmetric matrix is positive definite if and only if all its eigenvalues are non-negative.

### Exercise 3: Dot products are kernels



**Definition Dot Product** A dot product on a vector space  $\mathcal{H}$  is a symmetric bilinear form,

$$\begin{aligned} \langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} &\rightarrow \mathbb{R} \\ (\mathbf{x}, \mathbf{x}') &\mapsto \langle \mathbf{x}, \mathbf{x}' \rangle \end{aligned}$$

that is strictly positive definite; in other words, it has the property that for all  $\mathbf{x} \in \mathcal{H}$ ,  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$  with equality only for  $\mathbf{x} = 0$ .

Prove that dot products are positive definite kernels.

### Exercise 4: Kernel Logistic Regression



Consider a Logistic Regression (LR) model with the following loss function (cross entropy):

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N [y_i \log P_\theta(\hat{y} = 1 | \mathbf{x}_i) + (1 - y_i) \log (1 - P_\theta(\hat{y} = 1 | \mathbf{x}_i))] \quad (2)$$



## Exercise 7: Regression SVM



Consider the Lagrangian of the regression support vector machine (see [1] chapter 7.1.4 on SVMs for regression):

$$L(\mathbf{w}, b, \xi_i, \hat{\xi}_i, \alpha, \hat{\alpha}) = C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n (\beta_i \xi_i + \hat{\beta}_i \hat{\xi}_i) - \sum_{i=1}^n \alpha_i (\epsilon + \xi_i + y(\mathbf{x}_i) - y_i) - \sum_{i=1}^n \hat{\alpha}_i (\epsilon + \hat{\xi}_i - y(\mathbf{x}_i) + y_i), \quad (8)$$

(recall  $y(\mathbf{x}_i) = \mathbf{w}^\top \phi(\mathbf{x}_i) + b$ ), where  $E_\epsilon$  is the epsilon-insensitive error function:

$$E_\epsilon(y(\mathbf{x}) - y) = \begin{cases} 0 & \text{if } |y(\mathbf{x}) - y| < \epsilon \\ |y(\mathbf{x}) - y| - \epsilon & \text{otherwise} \end{cases} \quad (9)$$

with the largest accepted error  $\epsilon$ . We use Lagrange multipliers  $\alpha, \hat{\alpha}$  for the constraints with slack variables  $\xi_i, \hat{\xi}_i$ :

$$\begin{aligned} y_i &\leq y(\mathbf{x}_i) + \epsilon + \xi_i \\ y_i &\geq y(\mathbf{x}_i) - \epsilon - \hat{\xi}_i \end{aligned}$$

and  $\beta_i, \hat{\beta}_i$  to express the positivity constraints for  $\xi_i, \hat{\xi}_i$ .

By setting the derivatives of the Lagrangian with respect to  $\mathbf{w}, b, \xi_i$  and  $\hat{\xi}_i$  to zero and then back substituting to eliminate the corresponding variables, show that the dual Lagrangian is given by

$$\begin{aligned} \tilde{L}(\boldsymbol{\alpha}, \hat{\boldsymbol{\alpha}}) &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \hat{\alpha}_i) (\alpha_j - \hat{\alpha}_j) k(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad - \epsilon \sum_{i=1}^i (\alpha_i + \hat{\alpha}_i) + \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) y_i. \end{aligned} \quad (10)$$

with respect to  $\boldsymbol{\alpha}$  and  $\hat{\boldsymbol{\alpha}}$ . The kernel is defined as  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$ .

## References

- [1] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] B. Schölkopf and A. J. Smola. *Learning with kernels : support vector machines, regularization, optimization, and beyond*. Adaptive computation and machine learning. MIT Press, 2002.