

**STATISTICAL MACHINE LEARNING (WS2024)**  
**SUPPORT VECTOR MACHINES**

**Assignment 1.** Let  $\mathcal{T}^m = \{(\mathbf{x}^i, y^i) \in \mathbb{R}^n \times \{-1, +1\} \mid i = 1, \dots, m\}$  be a training set of linearly separable examples. In this case, the SVM algorithm finds parameters  $(\mathbf{w}^*, b^*) \in \mathbb{R}^{n+1}$  of a linear classifier

$$h(x; \mathbf{w}^*, b^*) = \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$$

by solving a convex quadratic program

$$(\mathbf{w}^*, b^*) \in \underset{(\mathbf{w}, b) \in \mathbb{R}^{n+1}}{\text{Arg min}} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \begin{cases} \langle \mathbf{w}, \mathbf{x}^i \rangle + b \geq 1, & i \in \mathcal{I}_+ \\ \langle \mathbf{w}, \mathbf{x}^i \rangle + b \leq -1, & i \in \mathcal{I}_- \end{cases} \quad (1)$$

where  $\mathcal{I}_+ = \{i \in \{1, \dots, m\} \mid y^i = +1\}$  and  $\mathcal{I}_- = \{i \in \{1, \dots, m\} \mid y^i = -1\}$  are indices of examples of the positive and the negative class, respectively.

**a)** Explain relation between the SVM for linearly separable examples, formulated in (1), and the generic SVM algorithm for learning from non-separable examples, which converts learning to

$$(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*) = \underset{\substack{(\mathbf{w}, b) \in \mathbb{R}^{n+1} \\ \boldsymbol{\xi} \in \mathbb{R}^m}}{\text{arg min}} \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \right)$$

subject to

$$\begin{cases} y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) \geq 1 - \xi_i, & i \in \{1, \dots, m\} \\ \xi_i \geq 0, & i \in \{1, \dots, m\} \end{cases}$$

**b)** For a given  $(\mathbf{w}, b) \in \mathbb{R}^{n+1}$  such that  $\mathbf{w} \neq 0$ , we can define so called margin

$$d(\mathbf{w}, b) = \min_{i \in \{1, \dots, m\}} \frac{y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b)}{\|\mathbf{w}\|}$$

which is a signed distance between the hyperplane  $\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$  and the closest examples in  $\mathcal{T}^m$ . Show that the hyperplane  $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^* = 0$  found by SVM algorithm (1) maximizes the margin, i.e. it holds that

$$d(\mathbf{w}^*, b^*) = \max_{\substack{\mathbf{w} \in \mathbb{R}^n \setminus \{0\} \\ b \in \mathbb{R}}} d(\mathbf{w}, b). \quad (2)$$

**c)** How can you compute the value of the maximal margin  $d(\mathbf{w}^*, b^*)$  from the solution  $(\mathbf{w}^*, b^*)$  ?

**Assignment 2.** Let the input observation be a vector  $\mathbf{x} \in \mathbb{R}^d$ . Let us consider a feature map  $\phi_q: \mathbb{R}^d \rightarrow \mathbb{R}^n$ ,  $n = d^q$ , whose entries are all possible  $q$ -th degree ordered products of the entries of  $\mathbf{x}$ . For example, if  $\mathbf{x} = (x_1, x_2, x_3)^T \in \mathbb{R}^3$  and  $q = 2$  then

$$\phi_q(\mathbf{x}) = \begin{pmatrix} x_1x_1 \\ x_2x_1 \\ x_3x_1 \\ x_1x_2 \\ x_2x_2 \\ x_3x_2 \\ x_1x_3 \\ x_2x_3 \\ x_3x_3 \end{pmatrix}$$

a) Show that for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  we can compute the dot product between  $\phi_q(\mathbf{x})$  and  $\phi_q(\mathbf{x}')$  as

$$\langle \phi_q(\mathbf{x}), \phi_q(\mathbf{x}') \rangle = \langle \mathbf{x}, \mathbf{x}' \rangle^q,$$

that is, as the dot product of the original vectors  $\mathbf{x}$  and  $\mathbf{x}'$  powered to  $q$ .

**Assignment 3.** Consider a feature map  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d(d+1)/2}$  whose entries are

$$\phi(\mathbf{x}) = \left( \begin{array}{cccccc} x_1^2, & \sqrt{2}x_1x_2, & \sqrt{2}x_1x_3, & \dots, & \sqrt{2}x_1x_d, \\ & x_2^2, & \sqrt{2}x_2x_3, & \dots, & \sqrt{2}x_2x_d, \\ & & & & \vdots \\ & & & & x_d^2 \end{array} \right)^T,$$

so that the features correspond to all possible products of unordered pairs of entries from  $\mathbf{x}$ , and the products of different entries are multiplied by a constant factor  $\sqrt{2}$ . For example, if  $\mathbf{x} = (x_1, x_2, x_3)^T \in \mathbb{R}^3$  then

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, x_2^2, \sqrt{2}x_2x_3, x_3^2)^T.$$

This feature map defines a kernel  $k(\mathbf{x}, \mathbf{x}') = \langle \phi'(\mathbf{x}), \phi'(\mathbf{x}') \rangle$  referred to as the homogeneous polynomial kernel of degree 2. Show that the kernel value equals to the square of the dot product of the input vectors, that is prove the identity

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi'(\mathbf{x}), \phi'(\mathbf{x}') \rangle = \langle \mathbf{x}, \mathbf{x}' \rangle^2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d.$$

**Assignment 4.** Assume we are given a training set of examples  $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \{+1, -1\}) \mid i = 1, \dots, m\}$  which is known to be linearly separable with respect to a feature map  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ . In this case, we can find parameters  $(\mathbf{w}, b) \in \mathbb{R}^{n+1}$  of a linear classifier  $h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$  which has zero training error by the Perceptron algorithm:

- (1)  $\mathbf{w} \leftarrow 0, b \leftarrow 0$
- (2) Find an example  $(x^u, y^u) \in \mathcal{T}^m$  whose label is incorrectly predicted by the current classifier, that is  $h(x^u; \mathbf{w}, b) \neq y^u$ .

- (3) If all examples are classified correctly exit the algorithm. Otherwise update the parameters by

$$\mathbf{w} \leftarrow \mathbf{w} + y^u \phi(x^u) \quad \text{and} \quad b \leftarrow b + y^u$$

and go to Step 2.

Assume that you cannot evaluate the feature map  $\phi(x)$  because it is either unknown or its evaluation is expensive. However, you know how to cheaply evaluate a kernel function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that  $k(x, x') = \langle \phi(x), \phi(x') \rangle, \forall x, x' \in \mathcal{X}$ . Show that you can still use the Perceptron algorithm to find a linear classifier with zero training error and that you can evaluate this classifier on any  $x \in \mathcal{X}$ .