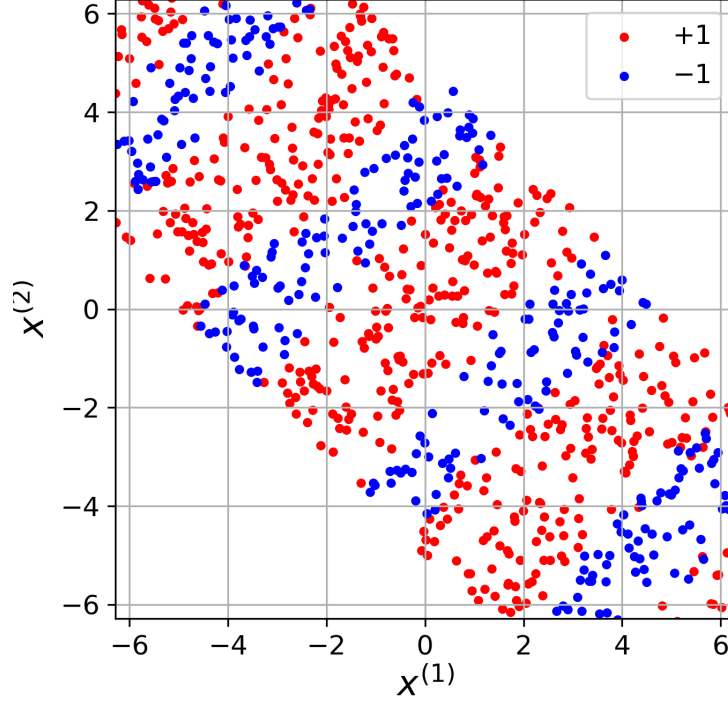


## Fourth assignment in Machine learning 1 – 2024 – Paper 1

### 1 Maximum marginal classifier (Deadline: 23rd October - MID DAY)



Assume a dataset  $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$  where  $\mathbf{x}_n \in \mathbb{R}^2$  and  $t_n \in \{-1, +1\}$ . Upon inspection of the data, we note a periodic pattern in the data (up to some exceptions). See the Figure above for an illustration.

To account for this periodicity, it seems reasonable to use some *trigonometric* functions to classify the data. In particular, we build classifier which will pick the class  $t_n = +1$  or  $-1$  depending on the sign of the following function:

$$y(\mathbf{x}_n) = w_0 \cos(x_n^{(1)} - x_n^{(2)}) + w_1 \sin(x_n^{(1)} - x_n^{(2)}) - \beta,$$

with  $\mathbf{w} = (w_0, w_1)^T \in \mathbb{R}^2$ ,  $\beta \in \mathbb{R}$  and  $\mathbf{x}_n = (x_n^{(1)}, x_n^{(2)})^T \in \mathbb{R}^2$ . To simplify our notation, we define the function  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  as

$$\phi(\mathbf{x}_n) = \begin{bmatrix} \cos(x_n^{(1)} - x_n^{(2)}) \\ \sin(x_n^{(1)} - x_n^{(2)}) \end{bmatrix},$$

such that  $y(\mathbf{x}_n) = \phi(\mathbf{x}_n)^T \mathbf{w} - \beta$ . Under the assumption that the two classes are perfectly separable, the distance between the decision boundary and any datapoint  $\mathbf{x}_n$  is given by:

$$\forall n, \quad \frac{t_n(\phi(\mathbf{x}_n)^T \mathbf{w} - \beta)}{\|\mathbf{w}\|_2} \geq 0. \quad (1)$$

As usual, it is convenient to leverage the scale invariance of this distance. Indeed,

$$\frac{\alpha}{\alpha} \cdot \frac{t_n(\phi(\mathbf{x}_n)^T \mathbf{w} - \beta)}{\|\mathbf{w}\|_2} = \frac{t_n(\phi(\mathbf{x}_n)^T (\alpha \mathbf{w}) - \alpha \beta)}{\|\alpha \mathbf{w}\|_2}.$$

Because of the arbitrary choice of  $\alpha$ , we are free to decide that we want to find a solution  $\mathbf{w}$  such that  $t_n(\phi(\mathbf{x}_n)^T \mathbf{w} - b) = 1$  for those points closest to the decision boundary, with  $b = \alpha\beta$ . In this case, the following constraint holds for all datapoints:

$$\forall n, \quad t_n(\phi(\mathbf{x}_n)^T \mathbf{w} - b) \geq 1.$$

Unfortunately, based on the Figure above, it does not seem that the data is perfectly separable, hence, we will introduce slack variables. We will now state the primal program that will find the optimal decision boundary:

$$\begin{aligned} \arg \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N \xi_n, \\ \text{s.t.} \quad & \forall n : t_n(\phi(\mathbf{x}_n)^T \mathbf{w} - b) \geq 1 - \xi_n, \quad (i) \\ & \forall n : \xi_n \geq 0. \quad (ii) \end{aligned}$$

Answer the following questions:

- (a) What does the dataset  $\mathcal{D}$  look like after applying transformation  $\phi$ ? You do not have to provide a graphical answer, a description with words suffices.  
*Hint:* Think of  $\phi$  as a composition of two transformations  $\phi = \phi_2 \circ \phi_1$  with  $\phi_1(\mathbf{x}) = [x^{(1)} - x^{(2)}, x^{(1)} + x^{(2)}]^T$  and  $\phi_2(\mathbf{x}) = [\cos(x^{(1)}), \sin(x^{(2)})]^T$ . [1 point]
- (b) How does the decision boundary  $y(\mathbf{x}) = 0$  look in the new, transformed space? How does the vector  $\mathbf{w}$  affect the decision boundary? What about  $b$ ? [1 point]
- (c) Introduce Lagrange multipliers for the constraints and write down the primal Lagrangian. You should use the target function  $\frac{1}{2} \|\mathbf{w}\|_2^2$  (note that you can answer this question even if you did not find a solution to the previous ones). Use the following notation:  $\{\lambda_n\}$  are the Lagrange multipliers for the first constraint (i) and  $\{\mu_n\}$  for the second (ii). [1 point]
- (d) How many KKT conditions are there? Write down all KKT conditions (not the *stationary conditions*). [2 points]
- (e) Derive the dual Lagrangian and specify the dual optimization problem. That is, derive the *stationary conditions* (by computing  $\partial \ell / \partial \rho = 0$ , where  $\ell$  is the primal Lagrangian and  $\rho$  is a primal variable) and use the resulting identities to eliminate the primal variables  $\{\mathbf{w}, b, \xi_1, \dots, \xi_N\}$ . Do not forget to specify the constraints on the remaining dual variables. [3 points]

- (f) Note that, because we have applied the transformation  $\phi$ , we have already written down a kernelized dual Lagrangian. This results in a dual Lagrangian that no longer depends on  $\mathbf{x}_n^T \mathbf{x}_m$  but on  $\kappa(\mathbf{x}_n, \mathbf{x}_m)$ . What is the explicit form of  $\kappa(\mathbf{x}_n, \mathbf{x}_m)$  in your final solution to the dual lagrangian? [1 point]
- (g) Assume we have solved the dual program and found optimal values for  $\lambda_n$  and  $b$ . Now we want to apply our maximum margin classifier, to a new test case  $\mathbf{x}^*$ . Describe how to classify the new datapoint  $\mathbf{x}^*$  in dual space. [1 point]
- (h) Use the KKT conditions to derive for which values of  $\lambda_n, \mu_n$  and  $\xi_n$ , a data point  $\mathbf{x}_n$  lies (i) outside the margin and is correctly classified, (ii) on the margin (i.e., is support vector), (iii) on the right side of the decision boundary within the margin, and (iv) on the wrong side of the decision boundary. [2 points]
- (i) Find the optimal values for the other dual variables  $\mu_n$  assuming we have found optimal  $\lambda_n$ . Then, solve for the primal variables  $\{\mathbf{w}, b, \xi_1 \dots, \xi_N\}$ . Note that you do not need to know the dual optimization program to solve this question. You only need the KKT conditions. *Hint*: Consider stationary conditions from question (e) when solving for optimal  $\mu_n$ . [3 points]
- (j) Now assume you obtain datasets with the following pattern. For each dataset, provide the function  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  such that two classes are linearly separable after the transformation. [3 points]

