

Datamining 2 - Jennifer Probst - HW 6

Ex 1 a) for this variant of SVM (transductive SVM) we formulate the Lagrangian:

$$L((w, b, \xi, \xi_+, \alpha, \alpha^*, \beta, \beta^*)) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + C_-^* \sum_{j: y_j^* = -1} \xi_j^* + C_+^* \sum_{j: y_j^* = 1} \xi_j^* \\ - \sum_{i=1}^n \alpha_i (y_i (\langle \vec{w}, \vec{x}_i \rangle + b) - 1 + \xi_i) - \sum_{j: y_j^* = -1} \beta_j \xi_j \\ - \sum_{j: y_j^* = -1} \alpha_j^* (y_j^* (\langle \vec{w}, \vec{x}_j^* \rangle + b) - 1 + \xi_j^*) - \sum_{j: y_j^* = -1} \beta_j^* \xi_j^* \\ - \sum_{j: y_j^* = 1} \alpha_j^* (y_j^* (\langle \vec{w}, \vec{x}_j^* \rangle + b) - 1 + \xi_j^*) - \sum_{j: y_j^* = 1} \beta_j^* \xi_j^*$$

we now maximize the Lagrangian with respect to w, b, ξ, ξ_+, ξ_- taking the first derivatives and setting them to zero.

$$\frac{\partial L}{\partial b} \stackrel{!}{=} 0 = - \sum_{i=1}^n \alpha_i y_i - \sum_{j: y_j^* = -1} \alpha_j^* y_j^* - \sum_{j: y_j^* = 1} \alpha_j^* y_j^* \quad (*) \text{ by setting the indices } j \text{ of the } k \text{ test points to indices } n+1 \text{ to } n+k \text{ we can derive a more compact form where the } \alpha_i \text{ with } i \in (n+1, n+k) \text{ are weights belonging to } y_i \text{ with } i \in (n+1, n+k) \text{ from the original test set.}$$

$$= \sum_{i=1}^n \alpha_i y_i + \sum_{i=n+1}^{n+k} \alpha_i y_i$$

$$0 = \sum_{i=1}^{n+k} \alpha_i y_i$$

$$\frac{dL}{d\vec{w}} \stackrel{!}{=} 0 = \vec{w} - \sum_{i=1}^n \alpha_i y_i \vec{x}_i - \sum_{j: y_j^* = -1} \alpha_j^* y_j^* \vec{x}_j^* - \sum_{j: y_j^* = 1} \alpha_j^* y_j^* \vec{x}_j^* \quad \text{simplified with same argument as above}$$

$$\vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i + \sum_{i=n+1}^{n+k} \alpha_i y_i \vec{x}_i$$

$$\vec{w} = \sum_{i=1}^{n+k} \alpha_i y_i \vec{x}_i$$

$$\frac{\partial L}{\partial \xi_i} \stackrel{!}{=} 0 = C - \alpha_i - \beta_i \Rightarrow \beta_i = C - \alpha_i \quad \text{with the constraint } \beta_i \geq 0$$

$$\frac{\partial L}{\partial \xi_j} \stackrel{!}{=} 0 = C_-^* - \alpha_j^* - \beta_j^* \Rightarrow \beta_j^* = C_-^* - \alpha_j^* \Rightarrow 0 \leq \alpha_i \leq C \quad \forall i \in \{1, \dots, n\}$$

$$(*) \Rightarrow 0 \leq \alpha_i \leq C_-^* \quad \forall i \in \{n+1, \dots, n+k\} \text{ with } y_i = -1$$

$$\frac{\partial L}{\partial \xi_{+j}} \stackrel{!}{=} 0 = C_+^* - \alpha_j^* - \beta_j^* \Rightarrow \beta_j^* = C_+^* - \alpha_j^* \Rightarrow 0 \leq \alpha_i \leq C_+^* \quad \forall i \in \{n+1, \dots, n+k\} \text{ with } y_i = 1$$

$$(*) \Rightarrow 0 \leq \alpha_i \leq C_+^* \quad \forall i \in \{n+1, \dots, n+k\} \text{ for } y_i = 1$$

substituting these findings in the Lagrangian we can see that it simplifies a lot: (already using the indexing (*))

$$\begin{aligned}
 \max_{\alpha} L(\alpha) &= + C \sum_{i=1}^{n+k} \xi_i - \sum_{i=1}^{n+k} \beta_i \xi_i - \sum_{i=1}^{n+k} d_i (y_i (\langle \vec{w}, \vec{x}_i \rangle + b) - 1 + \xi_i) + \frac{1}{2} \|\vec{w}\|^2 \\
 \left(\begin{array}{l} \beta_i = C - \alpha_i \\ \sum_{i=1}^{n+k} \alpha_i y_i = 0 \end{array} \right) &= + C \sum_{i=1}^{n+k} \xi_i - \sum_{i=1}^{n+k} (C - \alpha_i) \xi_i - \sum_{i=1}^{n+k} d_i y_i (\langle \vec{w}, \vec{x}_i \rangle + b) + \sum_{i=1}^{n+k} d_i - \sum_{i=1}^{n+k} \alpha_i \xi_i + \frac{1}{2} \|\vec{w}\|^2 \\
 \left(\begin{array}{l} \text{using:} \\ \vec{w} = \sum_{i=1}^{n+k} \alpha_i y_i \vec{x}_i \end{array} \right) &= \sum_{i=1}^{n+k} \alpha_i + \frac{1}{2} \left\langle \sum_{i=1}^{n+k} \alpha_i y_i \vec{x}_i, \sum_{j=1}^{n+k} \alpha_j y_j \vec{x}_j \right\rangle - \left\langle \sum_{i=1}^{n+k} \alpha_i y_i \vec{x}_i, \vec{w} \right\rangle \\
 &= \sum_{i=1}^{n+k} \alpha_i - \frac{1}{2} \sum_{i=1}^{n+k} \sum_{j=1}^{n+k} \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle
 \end{aligned}$$

→ The dual representation is:

$$\max_{\alpha} w(\alpha) = \sum_{i=1}^{n+k} \alpha_i - \frac{1}{2} \sum_{i=1}^{n+k} \sum_{j=1}^{n+k} \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

$$\begin{aligned}
 \text{s.t. } \sum_{i=1}^{n+k} \alpha_i y_i &= 0 \quad \text{with} \quad 0 \leq \alpha_i \leq C \quad \forall i \in \{1, \dots, n\} \\
 & \quad 0 \leq \alpha_i \leq C^* \quad \forall i \in \{n+1, \dots, n+k\} \text{ and } y_i = -1 \\
 & \quad 0 \leq \alpha_i \leq C^* \quad \forall i \in \{n+1, \dots, n+k\} \text{ and } y_i = 1
 \end{aligned}$$

2c) A Hilbert space H is:

- a vector space over real numbers \mathbb{R} which has a norm $\|\cdot\|_H$ and which is complete in that norm.
- In addition the norm comes from an inner product $\langle \cdot, \cdot \rangle_H$ (i.e. $\|f\|_H = \sqrt{\langle f, f \rangle_H}$)
- A kernel is a function $K: X \times X \rightarrow \mathbb{R}$ if there exists a Hilbert space H over \mathbb{R} and the map $\Phi: X \rightarrow H$ such that $\forall x, y \in X: K(x, y) = \langle \Phi(x), \Phi(y) \rangle_H$
- The evaluation functional at $x \in X$ in H is defined as the map $L_x: H \rightarrow \mathbb{R}$ such that for all functions $f \in H$ it holds: $L_x(f) = f(x)$
- we then call H a reproducing kernel Hilbert space if $\forall x \in X$, L_x is bounded ($\exists M$ constant and finite) such that $\|L_x(f)\| = |f(x)| \leq M \|f\|_H$
- This Reproducing kernel Hilbert space is associated with a kernel where evaluation at any x from every function in the space is reproduced by taking an inner product with a function determined by the kernel: kernel
- According to the Riesz representation theorem: $\forall x \in X \exists$ some function $K_x \in H$ such that $f(x) = L_x(f) = \langle f, K_x \rangle_H \quad \forall f \in H$.
- But as K_x is also a function on H , for any $y \in X: K_x(y) = L_y(K_x) = \langle K_x, K_y \rangle_H$
- if we now define $K: X \times X \rightarrow \mathbb{R}: \boxed{K(x, y) = \langle K_x, K_y \rangle_H}$
- with $\Phi(x) = K_x \rightarrow$ the kernel reproduces the inner product $K(x, y) = \langle \Phi(x), \Phi(y) \rangle_H$

2d) I would choose them by using cross validation for C_1 and C_2 .
 C_1 would probably be larger than C_2 , as the label for test data set is not known for certain → small C_2 . For the training data I would try a range of values for C_1 in the cross validation (only use part of ds to train & predict) to get a trade off b/w misclassification & margin size.