

Book Exercise 6.7: Verify the results 6.17 and 6.18 for constructing valid kernels.

(6.17) Given 2 valid kernels $k_1(x, x')$, $k_2(x, x')$ then

$K(x, x') = k_1(x, x') + k_2(x, x')$ is a valid kernel.

↳ let $k_1(x, x') = \phi_1(x)^T \phi_1(x')$ and $k_2(x, x') = \phi_2(x)^T \phi_2(x')$

We first show symmetry:

$$K(x, x') = k_1(x, x') + k_2(x, x') = k_1(x', x) + k_2(x', x) = K(x', x) \checkmark$$

Need Ψ such that $K(x, x') = \Psi(x)^T \Psi(x') = \phi_1(x)^T \phi_1(x') + \phi_2(x)^T \phi_2(x')$

So, we can take $\Psi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}$ and we see:

$$\Psi(x)^T \Psi(x') = \begin{bmatrix} \phi_1(x) & \phi_2(x) \end{bmatrix} \begin{bmatrix} \phi_1(x') \\ \phi_2(x') \end{bmatrix} = \phi_1(x)^T \phi_1(x') + \phi_2(x)^T \phi_2(x') = K(x, x') \checkmark$$

(6.18) $K(x, x') = k_1(x, x') k_2(x, x')$ is a valid kernel. Show it.

↳ let $k_1(x, x') = \phi(x)^T \phi(x')$, $k_2(x, x') = \phi'(x)^T \phi'(x')$. Then,

$$K(x, x') = [\phi(x)^T \phi(x')] [\phi'(x)^T \phi'(x')] \Rightarrow \sum_{i,j} [\phi_i(x) \phi_j(x')] [\phi'_i(x) \phi'_j(x')]$$

so we let $\Psi(x) = (\phi_i(x) \phi'_j(x))_{i,j}$ and we see

$$\Psi(x)^T \Psi(x') = K_1(x, x') K_2(x, x') = K(x, x') \checkmark$$

7.2:

Show that if the 1 on the right hand side of constraint 7.5 is replaced with $\gamma > 0$, the solution for the maximum margin hyperplane is unchanged.

⇒ Then, the constraint would be: $t_n(w^T \phi(x_n) + b) \geq \gamma$ $n=1, \dots, N$

If we look at the case where equality holds, there will always be one point such that the constraint is active. Once the margin is maximized, there will be at least two points with the active constraint.

The maximization problem is then:

$$\arg \max_{w, b} \left\{ \frac{1}{\|w\|} \min_n (t_n(w^T \phi(x_n) + b)) \right\} \text{ subject to } t_n(w^T \phi(x_n) + b) \geq \gamma$$

We can see then that the maximization problem boils down to maximizing $\|w\|^{-1}$, which is equivalent to minimizing $\|w\|^2$.

So we can formulate the problem in Lagrangian terms as:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n \{t_n (w^T \phi(x_n) + b) - \gamma\} \quad \text{where the } \frac{1}{2} \text{ is for convenience and } \alpha_n \text{ are the Lagr. coeff.}$$

① Take derivative wrt w :

$$\hookrightarrow \|w\| - \sum_{n=1}^N \alpha_n t_n \phi(x_n) = 0 \Rightarrow w = \sum_{n=1}^N \alpha_n t_n \phi(x_n) \quad \checkmark$$

② Take derivative wrt b :

$$\frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n t_n w^T \phi(x_n) + \alpha_n t_n b - \alpha_n \gamma \Rightarrow - \sum_{n=1}^N \alpha_n t_n = 0$$

$$\Rightarrow \sum_{n=1}^N \alpha_n t_n = 0 \quad \checkmark$$

Thus, we see the same result holds with $\gamma > 0$ instead of 1. ■

3. Given a dataset $D = \{(x_n, t_n)\}_{n=1}^N$ with $x_n \in \mathbb{R}^D$ and $t_n \in \{-1, 1\}$ for all n . The following is a formulation of soft-margin L_2 -SVM, a variant of the standard SVM obtained by squaring the hinge loss:

$$\begin{aligned} & x_n, t_n \quad n=1, \dots, N \quad t_n \in \{-1, 1\} \\ & \underset{w, b, \xi}{\text{minimize}} \quad \frac{1}{2} \|w\|_2^2 + C \sum_{n=1}^N \xi_n^2 \\ & \text{subj. to} \quad t_n (w^T \phi(x_n) + b) \geq 1 - \xi_n \quad \forall n \\ & \quad \quad \quad \xi_n \geq 0 \quad \forall n \end{aligned}$$

(a) Show that removing the last set of constraints $\{\xi_n \geq 0 \quad \forall n\}$ does not change the optimal solution to the problem above. Provide a complete proof.

(b) Describe the role of the hyperparameter $C \geq 0$. PS, 352

a) Pf: The last constraint $\xi_n \geq 0$ does not change the solution because of the definition/formulation of ξ_n .

$$\xi_n = \begin{cases} 0 & \text{if } x_n \text{ on correct side of margin} \\ |t_n - y(x_n)| & \text{otherwise} \end{cases}$$

Thus, we see that $\xi_n \geq 0 \quad \forall n$ implicitly, because:

$$\xi_n = 0 \geq 0$$

$$\text{OR } \xi_n = |t_n - y(x_n)| \geq 0$$

and these two cases cover all $n \in [1, N]$.

Therefore, the optimization problem

$$\text{minimize } \frac{1}{2} \|w\|_2^2 + C \sum_{n=1}^N \xi_n^2 \quad \text{subject to } \begin{matrix} t_n y(x_n) \geq 1 - \xi_n & \forall n \\ \xi_n \geq 0 & \forall n \end{matrix} \quad \text{and}$$

is redundant in terms of the last explicit constraint, and so it

can equivalently be rewritten as

$$\operatorname{argmin} \frac{1}{2} \|w\|_2^2 + C \sum_{n=1}^N \xi_n^2 \quad \text{subject to } t_n v(x_n) \geq 1 - \xi_n \quad \forall n$$

and it will have the same result, since $\xi_n \geq 0$ implicitly. ✓

- b) The hyperparameter $C \geq 0$ controls the trade-off between the slack variable penalty and the margin. The hyperparameter C is similar to a regularization parameter in that it controls the model complexity and overfitting, but it works inversely to a regularization parameter. It is inverse in the fact that a higher C value correlates to a less overfitted model. When $C = 0$ we have hard margin SVM, and when $C \rightarrow \infty$ we recover the linearly separable case of SVM.