

Kernel Based Support Vector Machines

Proof: Weights Lie in Space Spanned by $\phi(\mathbf{x}^j)$

©Barry Van Veen 2019

Background: Classifier training features and labels are $\mathbf{x}^i, d^i, i = 1, 2, \dots, N$. Classification in a high-dimensional feature space is performed using the mapping $\phi(\mathbf{x})$ as $\hat{d}(\mathbf{x}) = \text{sign} \{ \phi^T(\mathbf{x})\mathbf{w} \}$.

Claim: The weights \mathbf{w} that satisfy

$$\min_{\mathbf{w}} f(\mathbf{w}) = \min_{\mathbf{w}} \sum_{i=1}^N (1 - d^i \phi^T(\mathbf{x}^i)\mathbf{w})_+ + \lambda \|\mathbf{w}\|_2^2 \quad (1)$$

are of the form

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j$$

Proof: The proof proceeds by adding a component to \mathbf{w} that is orthogonal to the space spanned by the vectors $\{\phi(\mathbf{x}^j), j = 1, 2, \dots, N\}$ and then showing that component must be zero at the minimum of Eq. 1.

Suppose

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp$$

where ϕ^\perp is orthogonal to the space spanned by $\{\phi(\mathbf{x}^j), j = 1, 2, \dots, N\}$, that is, $\phi^T(\mathbf{x}^j)\phi^\perp = 0, j = 1, 2, \dots, N\}$. Note that any vector \mathbf{w} can be expressed as a sum of a component in the space spanned by the $\phi(\mathbf{x}^j)$ and a component orthogonal to that same space.

The optimization problem Eq. 1 may be rewritten

$$\begin{aligned} \min_{\mathbf{w}} f(\mathbf{w}) &= \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \phi^T(\mathbf{x}^i) \left(\sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right) \right)_+ + \lambda \left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right\|_2^2 \\ &= \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \left(\sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j + \phi^T(\mathbf{x}^i) \phi^\perp \right) \right)_+ \\ &\quad + \lambda \left(\sum_{i=1}^N \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_i \alpha_j + 2 \sum_{i=1}^N \alpha_i \phi^T(\mathbf{x}^i) \phi^\perp + \phi^{\perp T} \phi^\perp \right) \end{aligned} \quad (2)$$

$$\quad (3)$$

where in the second line we have used the identity

$$\left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right\|_2^2 = \left(\sum_{i=1}^N \phi(\mathbf{x}^i) \alpha_i + \phi^\perp \right)^T \left(\sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right)$$

and multiplied out the terms in the product.

Now use the fact that $\phi^T(\mathbf{x}^i) \phi^\perp = 0$ to reexpress the optimization problem as

$$\min_{\mathbf{w}} f(\mathbf{w}) = \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_i \alpha_j + \lambda \phi^{\perp T} \phi^\perp \quad (4)$$

$$= \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j \right)_+ + \lambda \left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j \right\|_2^2 + \lambda \|\phi^\perp\|_2^2 \quad (5)$$

The only term containing ϕ^\perp is the last one, $\lambda \|\phi^\perp\|_2^2$, which is nonnegative since $\lambda > 0$. Consequently, we conclude the minimum is attained when $\phi^\perp = \mathbf{0}$ and thus

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j$$