

Brief Article

The Author

October 9, 2013

I 6.1

Originally, we optimize $J(\mathbf{w}) = 1/2(\Phi\mathbf{w} - \mathbf{t})^T(\Phi\mathbf{w} - \mathbf{t})$. Now we set $\mathbf{w} = \Phi^T \mathbf{a}$ (since the optimal \mathbf{w} must take this form) and instead optimize \mathbf{a} . That is (6.7) in PRML pp. 293. The solution of (6.7) takes the form $K(K + \lambda I)\mathbf{a} = K\mathbf{t}$, and we can see the solution of \mathbf{a} is $(K + \lambda I)^{-1}(t + N(K)\lambda)$. Every one of them is fine. We can split \mathbf{a} into a part in the span of K and the other part in the null space of K , and we can check that by only preserving the span part, $J(\mathbf{a})$ is the same. So we have $\mathbf{a} = \text{vect}\Phi\mathbf{u}$, and see the solution manual. Actually, I think it's useless to talk about the rank of matrix $\Phi^T \Phi$. The point is, if we have a solution to *alpha*, we can use it to get a solution of $\mathbf{w} = \Phi^T \text{alpha}$.

2 6.2

Certain we can see that $\mathbf{w} = \sum_n a_n t_n \phi(\mathbf{x}_n)$, where a_i is the number of addition of this signed (normalized) sample. The learning rule is, whenever we classify $\phi(\mathbf{x}_n)$ wrongly, we make a_n bigger by one. In prediction, we use $\mathbf{w} = \Phi^T \boldsymbol{\alpha}$ to get $\mathbf{w}^T \phi(\mathbf{x}) = \boldsymbol{\alpha}^T \Phi \phi(\mathbf{x})$. Clearly the feature vector enters only in the form of the kernel function.