

9.2.3 Kernel Methods

We have all the tools together now to make an exciting step. Let us summarize our findings. We are interested in regularized estimation problems of the form (9.5) where $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ is linear, examples include the soft margin SVM and MAP for logistic regression. Here is a mad idea. Suppose we use a huge number of features p , maybe even infinitely many. Before figuring out how this could be done, let us first see whether this makes any sense in principle. After all, we have to store $\mathbf{w} \in \mathbb{R}^p$ and evaluate $\phi(\mathbf{x}) \in \mathbb{R}^p$. Do we? In the previous section, we learned that we can always represent $\mathbf{w} = \Phi^T \alpha$, where $\alpha \in \mathbb{R}^n$, and our dataset is finite. Moreover, the error function in (9.5) depends on

$$y = \Phi \mathbf{w} + b = \Phi \Phi^T \alpha + b$$

only, and $\Phi \Phi^T$ is just an $\mathbb{R}^{n \times n}$ matrix. Finally, the Tikhonov regularizer is given by

$$\frac{\nu}{2} \|\mathbf{w}\|^2 = \frac{\nu}{2} \|\Phi^T \alpha\|^2 = \frac{\nu}{2} \alpha^T \Phi \Phi^T \alpha,$$

it also only depends on $\Phi \Phi^T$. Finally, once we are done and found (α_*, b_*) , where $\mathbf{w}_* = \Phi^T \alpha_*$, we can predict on new inputs \mathbf{x} with

$$y^*(\mathbf{x}) = \mathbf{w}_*^T \phi(\mathbf{x}) + b_* = \alpha_*^T \Phi \phi(\mathbf{x}) + b_*.$$

We need finite quantities only in order to make our idea work, namely the matrix $\Phi \Phi^T$ during training, and the mapping $\Phi \phi(\mathbf{x})$ for predictions later on, to be evaluated at finitely many \mathbf{x} . This is the basic observation which makes *kernelization* work.

The entries of $\Phi \Phi^T$ are $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, while $[\Phi \phi(\mathbf{x})] = [\phi(\mathbf{x}_i)^T \phi(\mathbf{x})]$. We can write

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}'),$$

a *kernel function*. It is now clear that given the kernel function $K(\mathbf{x}, \mathbf{x}')$, we never need to access the underlying $\phi(\mathbf{x})$. In fact, we can forget about the dimensionality p and vectors of this size altogether. What makes $K(\mathbf{x}, \mathbf{x}')$ a kernel function? It must be the inner product in some feature space, but what does that imply? Let us work out some properties. First, a kernel function is obviously symmetric: $K(\mathbf{x}', \mathbf{x}) = K(\mathbf{x}, \mathbf{x}')$. Second, consider some arbitrary set $\{\mathbf{x}_i\}$ of n input points and construct the *kernel matrix* $\mathbf{K} = [K(\mathbf{x}_i, \mathbf{x}_j)] \in \mathbb{R}^{n \times n}$. Also, denote $\Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^T \in \mathbb{R}^{n \times p}$. Then,

$$\alpha^T \mathbf{K} \alpha = \alpha^T \Phi \Phi^T \alpha = \|\Phi^T \alpha\|^2 \geq 0.$$

In other words, the kernel matrix \mathbf{K} is positive semidefinite (see Section 6.3). This property defines kernel functions. $K(\mathbf{x}, \mathbf{x}')$ is a kernel function if the kernel matrix $\mathbf{K} = [K(\mathbf{x}_i, \mathbf{x}_j)]$ for any finite set of points $\{\mathbf{x}_i\}$ is symmetric positive semidefinite. An important subfamily are the *infinite-dimensional* or *positive definite kernel functions*. A member $K(\mathbf{x}, \mathbf{x}')$ of this subfamily is defined by all its kernel matrices $\mathbf{K} = [K(\mathbf{x}_i, \mathbf{x}_j)]$ being positive definite for any set $\{\mathbf{x}_i\}$ of any size. In particular, all kernel matrices are invertible. As we will see shortly, it is positive definite kernel functions which give rise to infinite-dimensional feature spaces, therefore to nonlinear kernel methods.

Problem: Can I use astronomically many features p ? How about $p = \infty$?

Approach: No problem! As long as you can efficiently compute the *kernel function* $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$, the representer theorem saves the day.

Hilbert Spaces and All That (*)

Before we give examples of kernel functions, a comment for meticulous readers (all others can safely skip this paragraph and move to the examples). How can we even talk about $\phi(\mathbf{x})^T \phi(\mathbf{x})$ if $p = \infty$? Even worse, what is $\Phi \in \mathbb{R}^{n \times p}$ in this case? In the best case, all this involves infinite sums, which may not converge. Rest assured that all this can be made rigorous within the framework of Hilbert function and functional spaces. In short, infinite dimensional vectors become functions, their transposes become functionals, and matrices become linear operators. A key result is Mercer's theorem for positive semidefinite kernel functions, which provides a construction for a feature map. However, with the exception of certain learning-theoretical questions, the importance of all this function space mathematics for down-to-earth machine learning is very limited. Historically, the point about the efforts of mathematicians like Hilbert, Schmidt and Riesz was to find conditions under which function spaces could be treated in the same simple way as finite-dimensional vector spaces, working out analogies for positive definite matrices, quadratic functions, eigendecomposition, and so on. Moreover, function spaces governing kernel methods are of the particularly simple reproducing kernel Hilbert type, where common pathologies like "delta functions" do not even arise. You may read about all that in [39] or other kernel literature, it will not play a role in this course. Just one warning which you will not find spelled out much in the SVM literature. The "geometry" in huge or infinite-dimensional spaces is dramatically different from anything we can draw or imagine. For example, in Mercer's construction of $K(\mathbf{x}, \mathbf{x}') = \sum_{j \geq 1} \phi_j(\mathbf{x}) \phi_j(\mathbf{x}')$, the different feature dimensions $j = 1, 2, \dots$ are by no means on equal terms, as far as concepts like distance or volume are concerned. For most commonly used infinite-dimensional kernel functions, the contributions $\phi_j(\mathbf{x}) \phi_j(\mathbf{x}')$ rapidly become extremely small, and only a small number of initial features determine most of the predictions. A good intuition about kernel methods is that they behave like (easy to use) linear methods of flexible dimensionality. As the number of data points n grows, a larger (but finite) number of the feature space dimensions will effectively be used.

Examples of Kernel Functions

Let us look at some examples. Maybe the simplest kernel function is $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$, the standard inner product. Moreover, for any finite-dimensional feature map $\phi(\mathbf{x}) \in \mathbb{R}^p$ ($p < \infty$), $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ is a kernel function. Since any kernel matrix of this type can at most have rank p , such kernel functions are positive semidefinite, but not positive definite. However, even for finite-dimensional kernels, it can be much simpler to work with $K(\mathbf{x}, \mathbf{x}')$ directly than to evaluate $\phi(\mathbf{x})$. For example, recall polynomial regression estimation from Section 4.1, giving rise to a polynomial feature map $\phi(x) = [1, x, \dots, x^r]^T$ for $x \in \mathbb{R}$. Now, if $\mathbf{x} \in \mathbb{R}^d$ is multivariate, a corresponding polynomial feature

map would consist of very many features. Is there a way around their explicit representation? Consider the *polynomial kernel*

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^r = \sum_{j_1, \dots, j_r} (x_{j_1} \dots x_{j_r}) (x'_{j_1} \dots x'_{j_r}).$$

For example, if $d = 3$ and $r = 2$, then

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= (x_1 x'_1 + x_2 x'_2 + x_3 x'_3)^2 = x_1^2 (x'_1)^2 + x_2^2 (x'_2)^2 + x_3^2 (x'_3)^2 \\ &\quad + 2(x_1 x_2)(x'_1 x'_2) + 2(x_1 x_3)(x'_1 x'_3) + 2(x_2 x_3)(x'_2 x'_3), \end{aligned}$$

a feature map of which is

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \sqrt{2}x_2x_3 \end{bmatrix}.$$

If $\mathbf{x} \in \mathbb{R}^d$, $K(\mathbf{x}, \mathbf{x}')$ is evaluated in $O(d)$, independent of r . Yet it is based on a feature map $\phi(\mathbf{x}) \in \mathbb{R}^{d^r}$, whose dimensionality⁶ scales exponentially in r . A variant is given by

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + \varepsilon)^r, \quad \varepsilon > 0,$$

which can be obtained by replacing \mathbf{x} by $[\mathbf{x}^T, \sqrt{\varepsilon}]^T$ above. The feature map now runs over all subranges $1 \leq j_1 \leq \dots \leq j_k \leq d$, $0 \leq k \leq r$.

A frequently used infinite-dimensional (positive definite) kernel is the *Gaussian* (or radial basis function, or squared exponential) kernel:

$$K(\mathbf{x}, \mathbf{x}') = e^{-\frac{\tau}{2} \|\mathbf{x} - \mathbf{x}'\|^2}, \quad \tau > 0. \quad (9.6)$$

We establish it as a kernel function in Section 9.2.4. The Gaussian is an example of a stationary kernel, these depend on $\mathbf{x} - \mathbf{x}'$ only. We can weight each dimension differently:

$$K(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2} \sum_{j=1}^d \tau_j (x_j - x'_j)^2}, \quad \tau_1, \dots, \tau_d > 0.$$

Free parameters in kernels are hyperparameters, much like C in the soft margin SVM or the noise variance σ^2 in Gaussian linear regression, choosing them is a model selection problem (Chapter 10).

Choosing the right kernel is much like choosing the right model. In order to do it well, you need to know your options. Kernels can be combined from others in many ways, [6, ch. 6.2] gives a good overview. It is also important to understand statistical properties implied by a kernel. For example, the Gaussian kernel produces extremely smooth solutions, while other kernels from the Matérn family are more flexible. Most books on kernel methods will provide some overview, see also [41].

One highly successful application domain for kernel methods concerns problems where input points \mathbf{x} have combinatorial structure, such as chains, trees, or

⁶More economically, we can run over all $1 \leq j_1 \leq \dots \leq j_r \leq d$.

graphs. Applications range from bioinformatics over computational chemistry to structured objects in computer vision. The rationale is that it is often simpler and much more computationally efficient to devise a kernel function $K(\mathbf{x}, \mathbf{x}')$ than a feature map $\phi(\mathbf{x})$. This field was seeded by independent work of David Haussler [22] and Chris Watkins.

A final remark concerns normalization. As noted in Chapter 2 and above in Section 9.1, it is often advantageous to use normalized feature maps $\|\phi(\mathbf{x})\| = 1$. What does this mean for a kernel?

$$K(\mathbf{x}, \mathbf{x}) = \phi(\mathbf{x})^T \phi(\mathbf{x}) = 1.$$

Therefore, a kernel function gives rise to a normalized feature map if its diagonal entries $K(\mathbf{x}, \mathbf{x})$ are all 1. For example, the Gaussian kernel (9.6) is normalized. Moreover, if $K(\mathbf{x}, \mathbf{x}')$ is a kernel, then so is

$$\frac{K(\mathbf{x}, \mathbf{x}')}{\sqrt{K(\mathbf{x}, \mathbf{x})K(\mathbf{x}', \mathbf{x}')}}.$$

(see Section 9.2.4), and the latter is normalized. It is a good idea to use normalized kernels in practice.

9.2.4 Techniques: Properties of Kernels (*)

In this section, we review a few properties of kernel functions and look at some more examples. The class of kernel functions has formidable closedness properties. If $K_1(\mathbf{x}, \mathbf{x}')$ and $K_2(\mathbf{x}, \mathbf{x}')$ are kernels, so are cK_1 for $c > 0$, $K_1 + K_2$ and $K_1 K_2$. You will have no problem confirming the first two. The third is shown at the end of this section. Moreover, $f(\mathbf{x})K_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$ is a kernel function as well, for any $f(\mathbf{x})$. This justifies kernel normalization, as discussed at the end of Section 9.2.3. If $K_r(\mathbf{x}, \mathbf{x}')$ is a sequence of kernel functions converging pointwise to $K(\mathbf{x}, \mathbf{x}') = \lim_{r \rightarrow \infty} K_r(\mathbf{x}, \mathbf{x}')$, then $K(\mathbf{x}, \mathbf{x}')$ is a kernel function as well. Finally, if $K(\mathbf{x}, \mathbf{x}')$ is a kernel and $\psi(\mathbf{y})$ is some mapping into \mathbb{R}^d , then $(\mathbf{y}, \mathbf{y}') \mapsto K(\psi(\mathbf{y}), \psi(\mathbf{y}'))$ is a kernel as well.

Let us show that the Gaussian kernel (9.6) is a valid kernel function. First, $(\mathbf{x}^T \mathbf{x}')^r$ is a kernel for every $r = 0, 1, 2, \dots$, namely the polynomial kernel from Section 9.2.3. By the way, $K(\mathbf{x}, \mathbf{x}') = 1$ is a kernel function, since its kernel matrices $\mathbf{1}\mathbf{1}^T$ are positive semidefinite. Therefore,

$$K_r(\mathbf{x}, \mathbf{x}') = \sum_{j=0}^r \frac{1}{j!} (\mathbf{x}^T \mathbf{x}')^j$$

are all kernels, and so is the limit $e^{\mathbf{x}^T \mathbf{x}'} = \lim_{r \rightarrow \infty} K_r(\mathbf{x}, \mathbf{x}')$. More general, if $K(\mathbf{x}, \mathbf{x}')$ is a kernel, so is $e^{K(\mathbf{x}, \mathbf{x}')}$. Now,

$$e^{-\frac{\tau}{2}\|\mathbf{x}-\mathbf{x}'\|^2} = e^{-\frac{\tau}{2}\|\mathbf{x}\|^2} e^{\tau\mathbf{x}^T \mathbf{x}'} e^{-\frac{\tau}{2}\|\mathbf{x}'\|^2}.$$

The middle is a kernel, and we apply our normalization rule with $f(\mathbf{x}) = e^{-\frac{\tau}{2}\|\mathbf{x}\|^2}$. The Gaussian kernel is infinite-dimensional (positive definite), although we will not show this here.

Another way to think about kernels is in terms of *covariance functions*. A random process is a set of random variables $a(\mathbf{x})$, one for each $\mathbf{x} \in \mathbb{R}^d$. Its covariance function is

$$K(\mathbf{x}, \mathbf{x}') = \text{Cov}[a(\mathbf{x}), a(\mathbf{x}')] = \mathbb{E}[(a(\mathbf{x}) - \mathbb{E}[a(\mathbf{x})])(a(\mathbf{x}') - \mathbb{E}[a(\mathbf{x}')])].$$

Covariance functions are kernel functions. For some set $\{\mathbf{x}_i\}$, let $\mathbf{a} = [a(\mathbf{x}_i) - \mathbb{E}[a(\mathbf{x}_i)]] \in \mathbb{R}^n$ be a random vector. Then, for any $\mathbf{v} \in \mathbb{R}^n$:

$$\mathbf{v}^T \mathbf{K} \mathbf{v} = \mathbf{v}^T \mathbb{E}[\mathbf{a} \mathbf{a}^T] \mathbf{v} = \mathbb{E}[(\mathbf{v}^T \mathbf{a})^2] \geq 0.$$

Finally, there are some symmetric functions which are not kernels. One example is

$$K(\mathbf{x}, \mathbf{x}') = \tanh(\alpha \mathbf{x}^T \mathbf{x}' + \beta).$$

In an attempt to make SVMs look like multi-layer perceptrons, this non-kernel was suggested and is shipped to this day in many SVM toolboxes⁷. Running the SVM with “kernels” like this spells trouble. Soft margin SVM is a convex optimization problem only if kernel matrices are positive semidefinite, codes will typically crash if that is not the case. A valid “neural networks” kernel is found in [47], derived from the covariance function perspective.

Finally, why is $K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$ a kernel? This argument is for interested readers only, it can be skipped at no loss. We have to show that for two positive semidefinite kernel matrices $\mathbf{K}_1, \mathbf{K}_2 \in \mathbb{R}^{n \times n}$ the Schur (or Hadamard) product $\mathbf{K}_1 \circ \mathbf{K}_2 = [K_1(\mathbf{x}_i, \mathbf{x}_j)K_2(\mathbf{x}_i, \mathbf{x}_j)]$ (Section 2.4.3) is positive semidefinite as well. To this end, we consider the Kronecker product $\mathbf{K}_1 \otimes \mathbf{K}_2 = [K_1(\mathbf{x}_i, \mathbf{x}_j)\mathbf{K}_2] \in \mathbb{R}^{n^2 \times n^2}$. This is positive semidefinite as well. Namely, we can write $\mathbf{K}_1 = \mathbf{V}_1 \mathbf{V}_1^T$, $\mathbf{K}_2 = \mathbf{V}_2 \mathbf{V}_2^T$, then

$$\mathbf{K}_1 \otimes \mathbf{K}_2 = (\mathbf{V}_1 \otimes \mathbf{V}_2)(\mathbf{V}_1 \otimes \mathbf{V}_2)^T.$$

But the Schur product is a square submatrix of $\mathbf{K}_1 \otimes \mathbf{K}_2$, a so called minor. In other words, for some index set $J \subset \{1, \dots, n^2\}$ of size $|J| = n$: $\mathbf{K}_1 \circ \mathbf{K}_2 = (\mathbf{K}_1 \otimes \mathbf{K}_2)_J$, so that for any $\mathbf{v} \in \mathbb{R}^n$:

$$\mathbf{v}^T (\mathbf{K}_1 \circ \mathbf{K}_2) \mathbf{v} = \mathbf{z}^T (\mathbf{K}_1 \otimes \mathbf{K}_2) \mathbf{z} \geq 0, \quad \mathbf{z} = \mathbf{I}_{\cdot, J} \mathbf{v} \in \mathbb{R}^{n^2}.$$

The same proof works to show that the positive definiteness of $\mathbf{K}_1, \mathbf{K}_2$ implies the positive definiteness of $\mathbf{K}_1 \circ \mathbf{K}_2$, a result due to Schur.

9.2.5 Summary

Let us summarize the salient points leading up to support vector machine binary classification. We started with the observation that for a linearly separable dataset, many different separating hyperplanes result in zero training error. Among all those potential solutions, which could arise as outcome of the perceptron algorithm, there is one which exhibits maximum stability against small displacements of patterns \mathbf{x}_i , by attaining the maximum margin $\gamma_{\mathcal{D}}(\mathbf{w}, b)$. Pursuing this lead, we addressed a number of problems:

⁷It is even on Wikipedia (en.wikipedia.org/wiki/Support_vector_machine).