

Kernels

The kernel function $k(\mathbf{x}, \mathbf{x}')$ is a real-valued function, and is often interpreted as a measure of similarity between the samples \mathbf{x} and \mathbf{x}' . In machine learning, kernels are often required to be symmetric and positive definite, and are typically employed when mapping non-linearly separable data in the input space into a higher-dimensional feature space where it can be linearly separated.

A key advantage of using kernel functions in machine learning is efficiency. Instead of explicitly defining a feature map and computing the inner product $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$, we can choose a kernel function that implicitly defines a feature map. When using a kernel function we operate in the original input space without having to compute the coordinates of the samples in the higher-dimension feature space. This is known as the kernel trick.

1 Positive Definite Kernels

Many machine learning methods require that the kernel function be symmetric and positive definite. This is because the corresponding Gram matrices (with which the algorithms work with in practice) are generally well-behaved, and possess some desirable properties: they are a requirement for Cholesky decomposition, they are invertible, have positive determinants, and they are associated with convex functions, making them very useful in optimisation algorithms.

A kernel function is positive definite if its associated Gram matrix, \mathbf{K} , with elements $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, is positive definite for any set of samples $\{\mathbf{x}_i\}_{i=1}^N$. Therefore, \mathbf{K} must satisfy

$$\mathbf{v}^T \mathbf{K} \mathbf{v} > 0 \tag{1}$$

for every non-zero real vector \mathbf{v} (see Appendix for example). If \mathbf{K} were positive semi-definite, it would have to satisfy

$$\mathbf{v}^T \mathbf{K} \mathbf{v} \geq 0 \tag{2}$$

The inequality (2) can alternately be expressed as

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j \mathbf{K}_{ij} > 0 \tag{3}$$

where $c_1, c_2, \dots, c_n \in \mathbb{R}$.

If the Gram matrix of a kernel is positive definite, Mercer's theorem states that it is possible to compute an eigenvector decomposition

$$\mathbf{K} = \mathbf{U}^T \mathbf{\Lambda} \mathbf{U} \quad (4)$$

$$= \mathbf{U}^T \mathbf{\Lambda}^{1/2} \mathbf{\Lambda}^{1/2} \mathbf{U} \quad (5)$$

$$= (\mathbf{\Lambda}^{1/2} \mathbf{U})^T (\mathbf{\Lambda}^{1/2} \mathbf{U}) \quad (6)$$

where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues $\lambda_i > 0$, and the rows of \mathbf{U} are the eigenvectors of \mathbf{K} .

An element of \mathbf{K} , can therefore be expressed as

$$\mathbf{K}_{ij} = (\mathbf{\Lambda}^{1/2} \mathbf{U}_{:i})^T (\mathbf{\Lambda}^{1/2} \mathbf{U}_{:j}) \quad (7)$$

If we define $\phi(\mathbf{x}_i) = \mathbf{\Lambda}^{1/2} \mathbf{U}_{:i}$, then we can write

$$\mathbf{K}_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (8)$$

Therefore, the elements of the Gram matrix \mathbf{K} are computed by performing an inner product of some feature vectors that are implicitly defined by the eigenvectors of \mathbf{K} . Generally, for the kernel function k , there exists a function ϕ , such that

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \quad (9)$$

where ϕ depends on the eigenfunctions of k , making the feature space potentially infinite-dimensional. ϕ is known as the feature map since it maps the sample \mathbf{x} in the input space to a higher-dimensional feature space.

Since inner products measures the similarity of vectors (if we normalise the vectors, the inner product will tell us the similarity of the directions of the vectors), kernel functions can be interpreted as a similarity measure between two elements of the input space.

2 The Polynomial Kernel

The polynomial kernel is defined as

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^d \quad (10)$$

where c is a constant, and d is the polynomial degree. Let's see how this works with an example:

Say we have a two-dimensional input space where

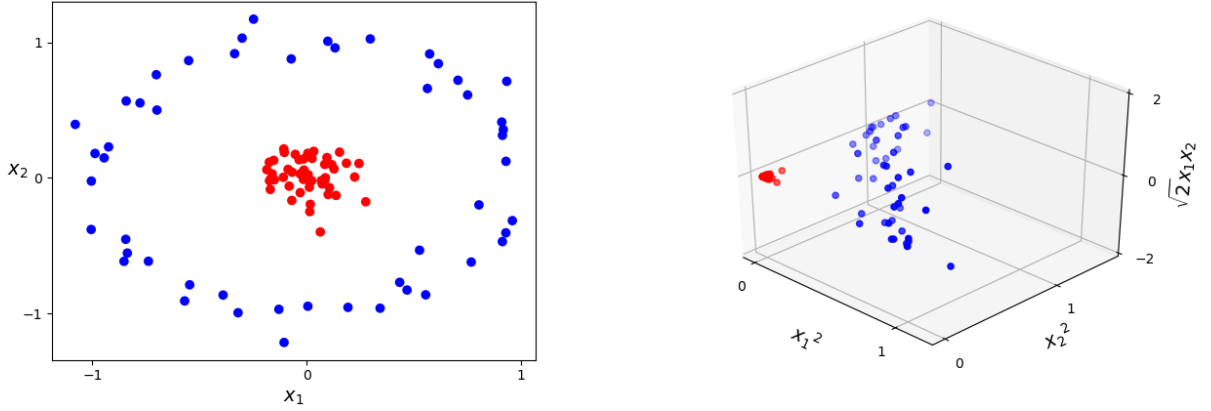


Figure 1: (left) A non-linearly separable, two-dimensional data set. (right) The data is mapped to a three-dimensional feature space, where it is linearly separable.

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \mathbf{x}' = \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix}$$

For a second-degree polynomial kernel we obtain

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= (x_1x'_1 + x_2x'_2 + c)^2 \\ &= (x_1^2x_1'^2 + x_2^2x_2'^2 + 2x_1x'_1x_2x'_2 + 2cx_1x'_1 + 2cx_2x'_2 + c^2) \end{aligned}$$

Therefore, from (9) we must have

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}cx_1 \\ \sqrt{2}cx_2 \\ c \end{pmatrix}$$

Hence, in this example we have mapped from a two-dimensional input space to a six-dimensional feature space. However, we do not actually need to compute the mapping explicitly, or even know it (as long as we know that it exists). The kernel function implicitly defines the mapping, and computing (10) will yield the same result as $\phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$. This is the kernel trick.

Now, consider the homogeneous case ($c = 0$). The second-degree polynomial kernel is

$$\begin{aligned}
k(\mathbf{x}, \mathbf{x}') &= (x_1x'_1 + x_2x'_2)^2 \\
&= (x_1^2x_1'^2 + x_2^2x_2'^2 + 2x_1x_1'x_2x_2')
\end{aligned}$$

Therefore, the feature map is

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$$

So, this time the mapping is to a three-dimensional feature space (see Figure 1).

3 The Radial Basis Function Kernel

The Radial Basis Function (RBF) kernel is defined as

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right) \quad (11)$$

where σ is a scaling parameter. The RBF kernel measures the Euclidean distance between the samples \mathbf{x} and \mathbf{x}' and returns a dot product in the feature space that depends on how close together the samples are. From equation (11) it is clear that the smaller the distance between the samples (the more similar the samples are to one another), the larger the value of the kernel function. The RBF kernel can be expanded as

$$\exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right) = \exp \left[-\frac{(\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')}{2\sigma^2} \right] \quad (12)$$

$$= \exp \left(\frac{\mathbf{x} \cdot \mathbf{x}'}{\sigma^2} \right) \exp \left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2} \right) \exp \left(-\frac{\|\mathbf{x}'\|^2}{2\sigma^2} \right) \quad (13)$$

$$= \left[1 + \frac{\mathbf{x} \cdot \mathbf{x}'}{\sigma^2} + \frac{1}{2!} \left(\frac{\mathbf{x} \cdot \mathbf{x}'}{\sigma^2} \right)^2 + \dots \right] \exp \left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2} \right) \exp \left(-\frac{\|\mathbf{x}'\|^2}{2\sigma^2} \right) \quad (14)$$

$$= \exp \left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2} \right) \exp \left(-\frac{\|\mathbf{x}'\|^2}{2\sigma^2} \right) \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{\mathbf{x} \cdot \mathbf{x}'}{\sigma^2} \right)^n \quad (15)$$

In the last line the two exponential terms are constant. The final term is just an infinite sum of the polynomial kernel (with $c = 0$) divided by some constants. Therefore, the mapping to this infinite-dimensional feature space can be expressed as the dot product of the feature maps, as is required.

If we define the homogeneous polynomial kernel of degree n as

$$\phi^n(\mathbf{x}) \cdot \phi^n(\mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}')^n \quad (16)$$

we can write the expanded RBF kernel as

$$\exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2}\right) \exp\left(-\frac{\|\mathbf{x}'\|^2}{2\sigma^2}\right) \sum_{n=0}^{\infty} \frac{\phi^n(\mathbf{x}) \cdot \phi^n(\mathbf{x}')}{n! \sigma^{2n}} \quad (17)$$

and the feature map as

$$\phi_{RBF}(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2}\right) \left(1, \frac{\phi^1(\mathbf{x})}{\sigma}, \frac{\phi^2(\mathbf{x})}{\sqrt{2!} \sigma^2}, \frac{\phi^3(\mathbf{x})}{\sqrt{3!} \sigma^3}, \dots\right)^T \quad (18)$$

For a more concrete example, let's say that \mathbf{x} and \mathbf{x}' are vectors in one-dimensional space with components x and x' , respectively. From equation (15) we obtain

$$\exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right) = \exp\left(-\frac{x^2}{2\sigma^2}\right) \exp\left(-\frac{x'^2}{2\sigma^2}\right) \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{xx'}{\sigma^2}\right)^n \quad (19)$$

$$= \exp\left(-\frac{x^2}{2\sigma^2}\right) \exp\left(-\frac{x'^2}{2\sigma^2}\right) \left(1 + \frac{xx'}{\sigma^2} + \frac{(xx')^2}{2! \sigma^4} + \frac{(xx')^3}{3! \sigma^6} + \dots\right) \quad (20)$$

Therefore, the feature map is

$$\phi_{RBF}(\mathbf{x}) = \exp\left(-\frac{x^2}{2\sigma^2}\right) \left(1, \frac{x}{\sigma}, \frac{x^2}{\sqrt{2!} \sigma^2}, \frac{x^3}{\sqrt{3!} \sigma^3}, \dots\right)^T \quad (21)$$

4 Support Vector Machines (SVMs) and Kernels

The dual form of the soft-margin SVM optimisation problem is

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^i \alpha^j y^i y^j \mathbf{x}^i \cdot \mathbf{x}^j - \sum_{i=1}^m \alpha^i \\ \text{subject to} \quad & \sum_{i=1}^m \alpha^i y^i = 0 \\ & 0 \leq \alpha^i \leq C, \quad i = 1, 2, \dots, m \end{aligned} \quad (22)$$

In the function to be minimised, the dot product occurs only between the samples \mathbf{x}^i and \mathbf{x}^j , and does not involve any of the other parameters. Therefore, if the data set is non-linearly separable, as well as employing the soft-margin approach, we can choose to map the samples

to a higher-dimensional space, with the feature map $\phi(\mathbf{x}^i)$ for the sample \mathbf{x}^i . Then, the function to be minimised (subject to the same constraints as above) would become

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^i \alpha^j y^i y^j \phi(\mathbf{x}^i) \cdot \phi(\mathbf{x}^j) - \sum_{i=1}^m \alpha^i \quad (23)$$

which is

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^i \alpha^j y^i y^j K(\mathbf{x}^i, \mathbf{x}^j) - \sum_{i=1}^m \alpha^i \quad (24)$$

and the weight vector solution to the soft-margin SVM problem would be

$$\mathbf{w} = \sum_{i=1}^m \alpha^i y^i \phi(\mathbf{x}^i) \quad (25)$$

References

Deisenroth, M. S., Faldó, A. A., Ong, C. S. (2020). *Mathematics for Machine Learning*. Cambridge University Press.

Mohri, M., Rostamizadeh, A., Talwalker, A. (2012). *Foundations of Machine Learning*. The MIT Press, Cambridge, Massachusetts.

Murphy, K. P., (2012). *Machine Learning: A Probabilistic Perspective*. The MIT Press, Cambridge, Massachusetts.

Shalev-Shwartz, S., Ben-David, S. (2014). *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, USA.

Appendix

A Simple Example of a Positive Definite Matrix

Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 2 & 6 \\ 6 & 20 \end{bmatrix} \tag{26}$$

The quadratic form $\mathbf{v}^T \mathbf{A} \mathbf{v}$ is expanded as

$$\mathbf{v}^T \mathbf{A} \mathbf{v} = [v_1, v_2] \begin{bmatrix} 2 & 6 \\ 6 & 20 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \tag{27}$$

$$= 2v_1^2 + 20v_2^2 + 12v_1v_2 \tag{28}$$

$$= 2(v_1 + 3v_2)^2 + v_2^2 \tag{29}$$

We can see that as long as $\mathbf{v} \neq \mathbf{0}$, $\mathbf{v}^T \mathbf{A} \mathbf{v} > 0$. Therefore, \mathbf{A} is positive definite.