

Before we can mine data, it is important to first find a suitable data representation that facilitates data analysis. For example, for complex data such as text, sequences, images, and so on, we must typically extract or construct a set of attributes or features, so that we can represent the data instances as multivariate vectors. That is, given a data instance \mathbf{x} (e.g., a sequence), we need to find a mapping ϕ , so that $\phi(\mathbf{x})$ is the vector representation of \mathbf{x} . Even when the input data is a numeric data matrix, if we wish to discover nonlinear relationships among the attributes, then a nonlinear mapping ϕ may be used, so that $\phi(\mathbf{x})$ represents a vector in the corresponding high-dimensional space comprising nonlinear attributes. We use the term *input space* to refer to the data space for the input data \mathbf{x} and *feature space* to refer to the space of mapped vectors $\phi(\mathbf{x})$. Thus, given a set of data objects or instances \mathbf{x}_i , and given a mapping function ϕ , we can transform them into feature vectors $\phi(\mathbf{x}_i)$, which then allows us to analyze complex data instances via numeric analysis methods.

Example 5.1 (Sequence-based Features). Consider a dataset of DNA sequences over the alphabet $\Sigma = \{A, C, G, T\}$. One simple feature space is to represent each sequence in terms of the probability distribution over symbols in Σ . That is, given a sequence \mathbf{x} with length $|\mathbf{x}| = m$, the mapping into feature space is given as

$$\phi(\mathbf{x}) = \{P(A), P(C), P(G), P(T)\}$$

where $P(s) = \frac{n_s}{m}$ is the probability of observing symbol $s \in \Sigma$, and n_s is the number of times s appears in sequence \mathbf{x} . Here the input space is the set of sequences Σ^* , and the feature space is \mathbb{R}^4 . For example, if $\mathbf{x} = ACAGCAGTA$, with $m = |\mathbf{x}| = 9$, since A occurs four times, C and G occur twice, and T occurs once, we have

$$\phi(\mathbf{x}) = (4/9, 2/9, 2/9, 1/9) = (0.44, 0.22, 0.22, 0.11)$$

Likewise, for another sequence $\mathbf{y} = AGCAAGCGAG$, we have

$$\phi(\mathbf{y}) = (4/10, 2/10, 4/10, 0) = (0.4, 0.2, 0.4, 0)$$

The mapping ϕ now allows one to compute statistics over the data sample to make inferences about the population. For example, we may compute the mean

symbol composition. We can also define the distance between any two sequences, for example,

$$\begin{aligned}\delta(\mathbf{x}, \mathbf{y}) &= \|\phi(\mathbf{x}) - \phi(\mathbf{y})\| \\ &= \sqrt{(0.44 - 0.4)^2 + (0.22 - 0.2)^2 + (0.22 - 0.4)^2 + (0.11 - 0)^2} = 0.22\end{aligned}$$

We can compute larger feature spaces by considering, for example, the probability distribution over all substrings or words of size up to k over the alphabet Σ , and so on.

Example 5.2 (Nonlinear Features). As an example of a nonlinear mapping consider the mapping ϕ that takes as input a vector $\mathbf{x} = (x_1, x_2)^T \in \mathbb{R}^2$ and maps it to a “quadratic” feature space via the nonlinear mapping

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^T \in \mathbb{R}^3$$

For example, the point $\mathbf{x} = (5.9, 3)^T$ is mapped to the vector

$$\phi(\mathbf{x}) = (5.9^2, 3^2, \sqrt{2} \cdot 5.9 \cdot 3)^T = (34.81, 9, 25.03)^T$$

The main benefit of this transformation is that we may apply well-known linear analysis methods in the feature space. However, because the features are nonlinear combinations of the original attributes, this allows us to mine nonlinear patterns and relationships.

Whereas mapping into feature space allows one to analyze the data via algebraic and probabilistic modeling, the resulting feature space is usually very high-dimensional; it may even be infinite dimensional. Thus, transforming all the input points into feature space can be very expensive, or even impossible. Because the dimensionality is high, we also run into the curse of dimensionality highlighted later in Chapter 6.

Kernel methods avoid explicitly transforming each point \mathbf{x} in the input space into the mapped point $\phi(\mathbf{x})$ in the feature space. Instead, the input objects are represented via their $n \times n$ pairwise similarity values. The similarity function, called a *kernel*, is chosen so that it represents a dot product in some high-dimensional feature space, yet it can be computed without directly constructing $\phi(\mathbf{x})$. Let \mathcal{I} denote the input space, which can comprise any arbitrary set of objects, and let $\mathbf{D} = \{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{I}$ be a dataset comprising n objects in the input space. We can represent the pairwise similarity values between points in \mathbf{D} via the $n \times n$ *kernel matrix*, defined as

$$\mathbf{K} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \cdots & K(\mathbf{x}_1, \mathbf{x}_n) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & K(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & K(\mathbf{x}_n, \mathbf{x}_2) & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

where $K: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$ is a *kernel function* on any two points in input space. However, we require that K corresponds to a dot product in some feature space. That is, for any

$\mathbf{x}_i, \mathbf{x}_j \in \mathcal{I}$, the kernel function should satisfy the condition

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (5.1)$$

where $\phi: \mathcal{I} \rightarrow \mathcal{F}$ is a mapping from the input space \mathcal{I} to the feature space \mathcal{F} . Intuitively, this means that we should be able to compute the value of the dot product using the original input representation \mathbf{x} , without having recourse to the mapping $\phi(\mathbf{x})$. Obviously, not just any arbitrary function can be used as a kernel; a valid kernel function must satisfy certain conditions so that Eq. (5.1) remains valid, as discussed in Section 5.1.

It is important to remark that the transpose operator for the dot product applies only when \mathcal{F} is a vector space. When \mathcal{F} is an abstract vector space with an inner product, the kernel is written as $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$. However, for convenience we use the transpose operator throughout this chapter; when \mathcal{F} is an inner product space it should be understood that

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \equiv \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

Example 5.3 (Linear and Quadratic Kernels). Consider the identity mapping, $\phi(\mathbf{x}) \rightarrow \mathbf{x}$. This naturally leads to the *linear kernel*, which is simply the dot product between two input vectors, and thus satisfies Eq. (5.1):

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

For example, consider the first five points from the two-dimensional Iris dataset shown in Figure 5.1a:

$$\mathbf{x}_1 = \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} \quad \mathbf{x}_2 = \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix} \quad \mathbf{x}_3 = \begin{pmatrix} 6.6 \\ 2.9 \end{pmatrix} \quad \mathbf{x}_4 = \begin{pmatrix} 4.6 \\ 3.2 \end{pmatrix} \quad \mathbf{x}_5 = \begin{pmatrix} 6 \\ 2.2 \end{pmatrix}$$

The kernel matrix for the linear kernel is shown in Figure 5.1b. For example,

$$K(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2 = 5.9 \times 6.9 + 3 \times 3.1 = 40.71 + 9.3 = 50.01$$

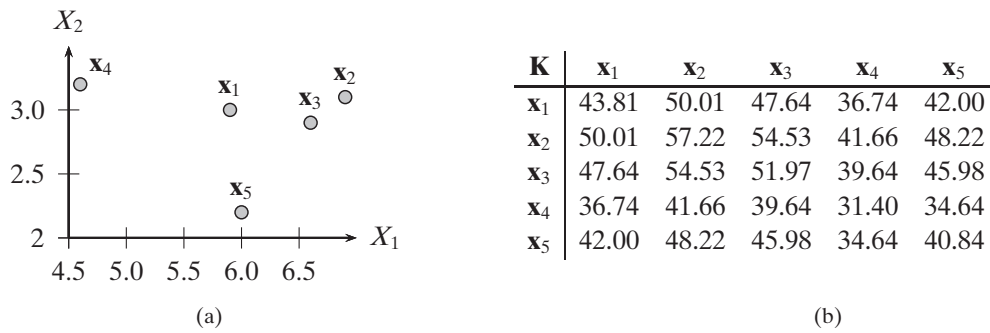


Figure 5.1. (a) Example points. (b) Linear kernel matrix.

Consider the quadratic mapping $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ from Example 5.2, that maps $\mathbf{x} = (x_1, x_2)^T$ as follows:

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^T$$

The dot product between the mapping for two input points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ is given as

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 x_2 y_2$$

We can rearrange the preceding to obtain the (homogeneous) *quadratic kernel* function as follows:

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{y}) &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 x_2 y_2 \\ &= (x_1 y_1 + x_2 y_2)^2 \\ &= (\mathbf{x}^T \mathbf{y})^2 \\ &= K(\mathbf{x}, \mathbf{y}) \end{aligned}$$

We can thus see that the dot product in feature space can be computed by evaluating the kernel in input space, without explicitly mapping the points into feature space. For example, we have

$$\begin{aligned} \phi(\mathbf{x}_1) &= (5.9^2, 3^2, \sqrt{2} \cdot 5.9 \cdot 3)^T = (34.81, 9, 25.03)^T \\ \phi(\mathbf{x}_2) &= (6.9^2, 3.1^2, \sqrt{2} \cdot 6.9 \cdot 3.1)^T = (47.61, 9.61, 30.25)^T \\ \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) &= 34.81 \times 47.61 + 9 \times 9.61 + 25.03 \times 30.25 = 2501 \end{aligned}$$

We can verify that the homogeneous quadratic kernel gives the same value

$$K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2)^2 = (50.01)^2 = 2501$$

We shall see that many data mining methods can be *kernelized*, that is, instead of mapping the input points into feature space, the data can be represented via the $n \times n$ kernel matrix \mathbf{K} , and all relevant analysis can be performed over \mathbf{K} . This is usually done via the so-called *kernel trick*, that is, show that the analysis task requires only dot products $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ in feature space, which can be replaced by the corresponding kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ that can be computed efficiently in input space. Once the kernel matrix has been computed, we no longer even need the input points \mathbf{x}_i , as all operations involving only dot products in the feature space can be performed over the $n \times n$ kernel matrix \mathbf{K} . An immediate consequence is that when the input data is the typical $n \times d$ numeric matrix \mathbf{D} and we employ the linear kernel, the results obtained by analyzing \mathbf{K} are equivalent to those obtained by analyzing \mathbf{D} (as long as only dot products are involved in the analysis). Of course, kernel methods allow much more flexibility, as we can just as easily perform non-linear analysis by employing nonlinear kernels, or we may analyze (non-numeric) complex objects without explicitly constructing the mapping $\phi(\mathbf{x})$.

Example 5.4. Consider the five points from Example 5.3 along with the linear kernel matrix shown in Figure 5.1. The mean of the five points in feature space is simply the mean in input space, as ϕ is the identity function for the linear kernel:

$$\boldsymbol{\mu}_\phi = \sum_{i=1}^5 \phi(\mathbf{x}_i) = \sum_{i=1}^5 \mathbf{x}_i = (6.00, 2.88)^T$$

Now consider the squared magnitude of the mean in feature space:

$$\|\boldsymbol{\mu}_\phi\|^2 = \boldsymbol{\mu}_\phi^T \boldsymbol{\mu}_\phi = (6.0^2 + 2.88^2) = 44.29$$

Because this involves only a dot product in feature space, the squared magnitude can be computed directly from \mathbf{K} . As we shall see later [see Eq. (5.12)] the squared norm of the mean vector in feature space is equivalent to the average value of the kernel matrix \mathbf{K} . For the kernel matrix in Figure 5.1b we have

$$\frac{1}{5^2} \sum_{i=1}^5 \sum_{j=1}^5 K(\mathbf{x}_i, \mathbf{x}_j) = \frac{1107.36}{25} = 44.29$$

which matches the $\|\boldsymbol{\mu}_\phi\|^2$ value computed earlier. This example illustrates that operations involving dot products in feature space can be cast as operations over the kernel matrix \mathbf{K} .

Kernel methods offer a radically different view of the data. Instead of thinking of the data as vectors in input or feature space, we consider only the kernel values between pairs of points. The kernel matrix can also be considered as a weighted adjacency matrix for the complete graph over the n input points, and consequently there is a strong connection between kernels and graph analysis, in particular algebraic graph theory.

5.1 KERNEL MATRIX

Let \mathcal{I} denote the input space, which can be any arbitrary set of data objects, and let $\mathbf{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathcal{I}$ denote a subset of n objects in the input space. Let $\phi: \mathcal{I} \rightarrow \mathcal{F}$ be a mapping from the input space into the feature space \mathcal{F} , which is endowed with a dot product and norm. Let $K: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$ be a function that maps pairs of input objects to their dot product value in feature space, that is, $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, and let \mathbf{K} be the $n \times n$ kernel matrix corresponding to the subset \mathbf{D} .

The function K is called a **positive semidefinite kernel** if and only if it is symmetric:

$$K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$$

and the corresponding kernel matrix \mathbf{K} for any subset $\mathbf{D} \subset \mathcal{I}$ is positive semidefinite, that is,

$$\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0, \text{ for all vectors } \mathbf{a} \in \mathbb{R}^n$$

which implies that

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \text{ for all } a_i \in \mathbb{R}, i \in [1, n] \quad (5.2)$$

We first verify that if $K(\mathbf{x}_i, \mathbf{x}_j)$ represents the dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ in some feature space, then K is a positive semidefinite kernel. Consider any dataset \mathbf{D} , and let $\mathbf{K} = \{K(\mathbf{x}_i, \mathbf{x}_j)\}$ be the corresponding kernel matrix. First, K is symmetric since the dot product is symmetric, which also implies that \mathbf{K} is symmetric. Second, \mathbf{K} is positive semidefinite because

$$\begin{aligned} \mathbf{a}^T \mathbf{K} \mathbf{a} &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j K(\mathbf{x}_i, \mathbf{x}_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &= \left(\sum_{i=1}^n a_i \phi(\mathbf{x}_i) \right)^T \left(\sum_{j=1}^n a_j \phi(\mathbf{x}_j) \right) \\ &= \left\| \sum_{i=1}^n a_i \phi(\mathbf{x}_i) \right\|^2 \geq 0 \end{aligned}$$

Thus, K is a positive semidefinite kernel.

We now show that if we are given a positive semidefinite kernel $K: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$, then it corresponds to a dot product in some feature space \mathcal{F} .

5.1.1 Reproducing Kernel Map

For the reproducing kernel map ϕ , we map each point $\mathbf{x} \in \mathcal{I}$ into a function in a *functional space* $\{f: \mathcal{I} \rightarrow \mathbb{R}\}$ comprising functions that map points in \mathcal{I} into \mathbb{R} . Algebraically this space of functions is an abstract vector space where each point happens to be a function. In particular, any $\mathbf{x} \in \mathcal{I}$ in the input space is mapped to the following function:

$$\phi(\mathbf{x}) = K(\mathbf{x}, \cdot)$$

where the \cdot stands for any argument in \mathcal{I} . That is, each object \mathbf{x} in the input space gets mapped to a *feature point* $\phi(\mathbf{x})$, which is in fact a function $K(\mathbf{x}, \cdot)$ that represents its similarity to all other points in the input space \mathcal{I} .

Let \mathcal{F} be the set of all functions or points that can be obtained as a linear combination of any subset of feature points, defined as

$$\begin{aligned} \mathcal{F} &= \text{span}\{K(\mathbf{x}, \cdot) \mid \mathbf{x} \in \mathcal{I}\} \\ &= \left\{ \mathbf{f} = f(\cdot) = \sum_{i=1}^m \alpha_i K(\mathbf{x}_i, \cdot) \mid m \in \mathbb{N}, \alpha_i \in \mathbb{R}, \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathcal{I} \right\} \end{aligned}$$

We use the dual notation \mathbf{f} and $f(\cdot)$ interchangeably to emphasize the fact that each point \mathbf{f} in the feature space is in fact a function $f(\cdot)$. Note that by definition the feature point $\phi(\mathbf{x}) = K(\mathbf{x}, \cdot)$ belongs to \mathcal{F} .

Let $\mathbf{f}, \mathbf{g} \in \mathcal{F}$ be any two points in feature space:

$$\mathbf{f} = f(\cdot) = \sum_{i=1}^{m_a} \alpha_i K(\mathbf{x}_i, \cdot) \quad \mathbf{g} = g(\cdot) = \sum_{j=1}^{m_b} \beta_j K(\mathbf{x}_j, \cdot)$$

Define the dot product between two points as

$$\mathbf{f}^T \mathbf{g} = f(\cdot)^T g(\cdot) = \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \alpha_i \beta_j K(\mathbf{x}_i, \mathbf{x}_j) \quad (5.3)$$

We emphasize that the notation $\mathbf{f}^T \mathbf{g}$ is only a convenience; it denotes the inner product $\langle \mathbf{f}, \mathbf{g} \rangle$ because \mathcal{F} is an abstract vector space, with an inner product as defined above.

We can verify that the dot product is *bilinear*, that is, linear in both arguments, because

$$\mathbf{f}^T \mathbf{g} = \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \alpha_i \beta_j K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{i=1}^{m_a} \alpha_i g(\mathbf{x}_i) = \sum_{j=1}^{m_b} \beta_j f(\mathbf{x}_j)$$

The fact that K is positive semidefinite implies that

$$\|\mathbf{f}\|^2 = \mathbf{f}^T \mathbf{f} = \sum_{i=1}^{m_a} \sum_{j=1}^{m_a} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

Thus, the space \mathcal{F} is a *pre-Hilbert space*, defined as a normed inner product space, because it is endowed with a symmetric bilinear dot product and a norm. By adding the limit points of all Cauchy sequences that are convergent, \mathcal{F} can be turned into a *Hilbert space*, defined as a normed inner product space that is complete. However, showing this is beyond the scope of this chapter.

The space \mathcal{F} has the so-called *reproducing property*, that is, we can evaluate a function $f(\cdot) = \mathbf{f}$ at a point $\mathbf{x} \in \mathcal{I}$ by taking the dot product of \mathbf{f} with $\phi(\mathbf{x})$, that is,

$$\mathbf{f}^T \phi(\mathbf{x}) = f(\cdot)^T K(\mathbf{x}, \cdot) = \sum_{i=1}^{m_a} \alpha_i K(\mathbf{x}_i, \mathbf{x}) = f(\mathbf{x})$$

For this reason, the space \mathcal{F} is also called a *reproducing kernel Hilbert space*.

All we have to do now is to show that $K(\mathbf{x}_i, \mathbf{x}_j)$ corresponds to a dot product in the feature space \mathcal{F} . This is indeed the case, because using Eq. (5.3) for any two feature points $\phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \in \mathcal{F}$ their dot product is given as

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \cdot)^T K(\mathbf{x}_j, \cdot) = K(\mathbf{x}_i, \mathbf{x}_j)$$

The reproducing kernel map shows that any positive semidefinite kernel corresponds to a dot product in some feature space. This means we can apply well known algebraic and geometric methods to understand and analyze the data in these spaces.

Empirical Kernel Map

The reproducing kernel map ϕ maps the input space into a potentially infinite dimensional feature space. However, given a dataset $\mathbf{D} = \{\mathbf{x}_i\}_{i=1}^n$, we can obtain a finite

dimensional mapping by evaluating the kernel only on points in \mathbf{D} . That is, define the map ϕ as follows:

$$\phi(\mathbf{x}) = \left(K(\mathbf{x}_1, \mathbf{x}), K(\mathbf{x}_2, \mathbf{x}), \dots, K(\mathbf{x}_n, \mathbf{x}) \right)^T \in \mathbb{R}^n$$

which maps each point $\mathbf{x} \in \mathcal{I}$ to the n -dimensional vector comprising the kernel values of \mathbf{x} with each of the objects $\mathbf{x}_i \in \mathbf{D}$. We can define the dot product in feature space as

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \sum_{k=1}^n K(\mathbf{x}_k, \mathbf{x}_i) K(\mathbf{x}_k, \mathbf{x}_j) = \mathbf{K}_i^T \mathbf{K}_j \quad (5.4)$$

where \mathbf{K}_i denotes the i th column of \mathbf{K} , which is also the same as the i th row of \mathbf{K} (considered as a column vector), as \mathbf{K} is symmetric. However, for ϕ to be a valid map, we require that $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$, which is clearly not satisfied by Eq. (5.4). One solution is to replace $\mathbf{K}_i^T \mathbf{K}_j$ in Eq. (5.4) with $\mathbf{K}_i^T \mathbf{A} \mathbf{K}_j$ for some positive semidefinite matrix \mathbf{A} such that

$$\mathbf{K}_i^T \mathbf{A} \mathbf{K}_j = K(\mathbf{x}_i, \mathbf{x}_j)$$

If we can find such an \mathbf{A} , it would imply that over all pairs of mapped points we have

$$\left\{ \mathbf{K}_i^T \mathbf{A} \mathbf{K}_j \right\}_{i,j=1}^n = \left\{ K(\mathbf{x}_i, \mathbf{x}_j) \right\}_{i,j=1}^n$$

which can be written compactly as

$$\mathbf{K} \mathbf{A} \mathbf{K} = \mathbf{K}$$

This immediately suggests that we take $\mathbf{A} = \mathbf{K}^{-1}$, the (pseudo) inverse of the kernel matrix \mathbf{K} . The modified map ϕ , called the *empirical kernel map*, is then defined as

$$\phi(\mathbf{x}) = \mathbf{K}^{-1/2} \cdot \left(K(\mathbf{x}_1, \mathbf{x}), K(\mathbf{x}_2, \mathbf{x}), \dots, K(\mathbf{x}_n, \mathbf{x}) \right)^T \in \mathbb{R}^n$$

so that the dot product yields

$$\begin{aligned} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) &= \left(\mathbf{K}^{-1/2} \mathbf{K}_i \right)^T \left(\mathbf{K}^{-1/2} \mathbf{K}_j \right) \\ &= \mathbf{K}_i^T \left(\mathbf{K}^{-1/2} \mathbf{K}^{-1/2} \right) \mathbf{K}_j \\ &= \mathbf{K}_i^T \mathbf{K}^{-1} \mathbf{K}_j \end{aligned}$$

Over all pairs of mapped points, we have

$$\left\{ \mathbf{K}_i^T \mathbf{K}^{-1} \mathbf{K}_j \right\}_{i,j=1}^n = \mathbf{K} \mathbf{K}^{-1} \mathbf{K} = \mathbf{K}$$

as desired. However, it is important to note that this empirical feature representation is valid only for the n points in \mathbf{D} . If points are added to or removed from \mathbf{D} , the kernel map will have to be updated for all points.

5.1.2 Mercer Kernel Map

In general different feature spaces can be constructed for the same kernel K . We now describe how to construct the Mercer map.

Data-specific Kernel Map

The Mercer kernel map is best understood starting from the kernel matrix for the dataset \mathbf{D} in input space. Because \mathbf{K} is a symmetric positive semidefinite matrix, it has real and non-negative eigenvalues, and it can be decomposed as follows:

$$\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

where \mathbf{U} is the orthonormal matrix of eigenvectors $\mathbf{u}_i = (u_{i1}, u_{i2}, \dots, u_{in})^T \in \mathbb{R}^n$ (for $i = 1, \dots, n$), and $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues, with both arranged in non-increasing order of the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$:

$$\mathbf{U} = \begin{pmatrix} | & | & \cdots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & \cdots & | \end{pmatrix} \quad \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{pmatrix}$$

The kernel matrix \mathbf{K} can therefore be rewritten as the spectral sum

$$\mathbf{K} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \cdots + \lambda_n \mathbf{u}_n \mathbf{u}_n^T$$

In particular the kernel function between \mathbf{x}_i and \mathbf{x}_j is given as

$$\begin{aligned} \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) &= \lambda_1 u_{1i} u_{1j} + \lambda_2 u_{2i} u_{2j} \cdots + \lambda_n u_{ni} u_{nj} \\ &= \sum_{k=1}^n \lambda_k u_{ki} u_{kj} \end{aligned} \quad (5.5)$$

where u_{ki} denotes the i th component of eigenvector \mathbf{u}_k . It follows that if we define the Mercer map ϕ as follows:

$$\phi(\mathbf{x}_i) = \left(\sqrt{\lambda_1} u_{1i}, \sqrt{\lambda_2} u_{2i}, \dots, \sqrt{\lambda_n} u_{ni} \right)^T \quad (5.6)$$

then $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$ is a dot product in feature space between the mapped points $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ because

$$\begin{aligned} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) &= \left(\sqrt{\lambda_1} u_{1i}, \dots, \sqrt{\lambda_n} u_{ni} \right) \left(\sqrt{\lambda_1} u_{1j}, \dots, \sqrt{\lambda_n} u_{nj} \right)^T \\ &= \lambda_1 u_{1i} u_{1j} + \cdots + \lambda_n u_{ni} u_{nj} = \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Noting that $\mathbf{U}_i = (u_{1i}, u_{2i}, \dots, u_{ni})^T$ is the i th row of \mathbf{U} , we can rewrite the Mercer map ϕ as

$$\phi(\mathbf{x}_i) = \sqrt{\mathbf{\Lambda}} \mathbf{U}_i \quad (5.7)$$

Thus, the kernel value is simply the dot product between scaled rows of \mathbf{U} :

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \left(\sqrt{\mathbf{\Lambda}} \mathbf{U}_i \right)^T \left(\sqrt{\mathbf{\Lambda}} \mathbf{U}_j \right) = \mathbf{U}_i^T \mathbf{\Lambda} \mathbf{U}_j$$

The Mercer map, defined equivalently in Eqs. (5.6) and (5.7), is obviously restricted to the input dataset \mathbf{D} , just like the empirical kernel map, and is therefore called the *data-specific Mercer kernel map*. It defines a data-specific feature space of dimensionality at most n , comprising the eigenvectors of \mathbf{K} .

Example 5.5. Let the input dataset comprise the five points shown in Figure 5.1a, and let the corresponding kernel matrix be as shown in Figure 5.1b. Computing the eigen-decomposition of \mathbf{K} , we obtain $\lambda_1 = 223.95$, $\lambda_2 = 1.29$, and $\lambda_3 = \lambda_4 = \lambda_5 = 0$. The effective dimensionality of the feature space is 2, comprising the eigenvectors \mathbf{u}_1 and \mathbf{u}_2 . Thus, the matrix \mathbf{U} is given as follows:

$$\mathbf{U} = \begin{pmatrix} & \mathbf{u}_1 & \mathbf{u}_2 \\ \mathbf{U}_1 & -0.442 & 0.163 \\ \mathbf{U}_2 & -0.505 & -0.134 \\ \mathbf{U}_3 & -0.482 & -0.181 \\ \mathbf{U}_4 & -0.369 & 0.813 \\ \mathbf{U}_5 & -0.425 & -0.512 \end{pmatrix}$$

and we have

$$\mathbf{\Lambda} = \begin{pmatrix} 223.95 & 0 \\ 0 & 1.29 \end{pmatrix} \quad \sqrt{\mathbf{\Lambda}} = \begin{pmatrix} \sqrt{223.95} & 0 \\ 0 & \sqrt{1.29} \end{pmatrix} = \begin{pmatrix} 14.965 & 0 \\ 0 & 1.135 \end{pmatrix}$$

The kernel map is specified via Eq. (5.7). For example, for $\mathbf{x}_1 = (5.9, 3)^T$ and $\mathbf{x}_2 = (6.9, 3.1)^T$ we have

$$\begin{aligned} \phi(\mathbf{x}_1) &= \sqrt{\mathbf{\Lambda}} \mathbf{U}_1 = \begin{pmatrix} 14.965 & 0 \\ 0 & 1.135 \end{pmatrix} \begin{pmatrix} -0.442 \\ 0.163 \end{pmatrix} = \begin{pmatrix} -6.616 \\ 0.185 \end{pmatrix} \\ \phi(\mathbf{x}_2) &= \sqrt{\mathbf{\Lambda}} \mathbf{U}_2 = \begin{pmatrix} 14.965 & 0 \\ 0 & 1.135 \end{pmatrix} \begin{pmatrix} -0.505 \\ -0.134 \end{pmatrix} = \begin{pmatrix} -7.563 \\ -0.153 \end{pmatrix} \end{aligned}$$

Their dot product is given as

$$\begin{aligned} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) &= 6.616 \times 7.563 - 0.185 \times 0.153 \\ &= 50.038 - 0.028 = 50.01 \end{aligned}$$

which matches the kernel value $K(\mathbf{x}_1, \mathbf{x}_2)$ in Figure 5.1b.

Mercer Kernel Map

For compact continuous spaces, analogous to the discrete case in Eq. (5.5), the kernel value between any two points can be written as the infinite spectral decomposition

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^{\infty} \lambda_k \mathbf{u}_k(\mathbf{x}_i) \mathbf{u}_k(\mathbf{x}_j)$$

where $\{\lambda_1, \lambda_2, \dots\}$ is the infinite set of eigenvalues, and $\{\mathbf{u}_1(\cdot), \mathbf{u}_2(\cdot), \dots\}$ is the corresponding set of orthogonal and normalized *eigenfunctions*, that is, each function $\mathbf{u}_i(\cdot)$ is a solution to the integral equation

$$\int K(\mathbf{x}, \mathbf{y}) \mathbf{u}_i(\mathbf{y}) d\mathbf{y} = \lambda_i \mathbf{u}_i(\mathbf{x})$$

and K is a continuous positive semidefinite kernel, that is, for all functions $a(\cdot)$ with a finite square integral (i.e., $\int a(\mathbf{x})^2 d\mathbf{x} < \infty$) K satisfies the condition

$$\iint K(\mathbf{x}_1, \mathbf{x}_2) a(\mathbf{x}_1) a(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2 \geq 0$$

We can see that this positive semidefinite kernel for compact continuous spaces is analogous to the discrete kernel in Eq. (5.2). Further, similarly to the data-specific Mercer map [Eq. (5.6)], the general Mercer kernel map is given as

$$\phi(\mathbf{x}_i) = \left(\sqrt{\lambda_1} \mathbf{u}_1(\mathbf{x}_i), \sqrt{\lambda_2} \mathbf{u}_2(\mathbf{x}_i), \dots \right)^T$$

with the kernel value being equivalent to the dot product between two mapped points:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

5.2 VECTOR KERNELS

We now consider two of the most commonly used vector kernels in practice. Kernels that map an (input) vector space into another (feature) vector space are called *vector kernels*. For multivariate input data, the input vector space will be the d -dimensional real space \mathbb{R}^d . Let \mathbf{D} comprise n input points $\mathbf{x}_i \in \mathbb{R}^d$, for $i = 1, 2, \dots, n$. Commonly used (nonlinear) kernel functions over vector data include the polynomial and Gaussian kernels, as described next.

Polynomial Kernel

Polynomial kernels are of two types: homogeneous or inhomogeneous. Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. The *homogeneous polynomial kernel* is defined as

$$K_q(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y}) = (\mathbf{x}^T \mathbf{y})^q \quad (5.8)$$

where q is the degree of the polynomial. This kernel corresponds to a feature space spanned by all products of exactly q attributes.

The most typical cases are the *linear* (with $q = 1$) and *quadratic* (with $q = 2$) kernels, given as

$$\begin{aligned} K_1(\mathbf{x}, \mathbf{y}) &= \mathbf{x}^T \mathbf{y} \\ K_2(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y})^2 \end{aligned}$$

The *inhomogeneous polynomial kernel* is defined as

$$K_q(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^q \quad (5.9)$$

where q is the degree of the polynomial, and $c \geq 0$ is some constant. When $c = 0$ we obtain the homogeneous kernel. When $c > 0$, this kernel corresponds to the feature space spanned by all products of at most q attributes. This can be seen from the binomial expansion

$$K_q(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^q = \sum_{k=1}^q \binom{q}{k} c^{q-k} (\mathbf{x}^T \mathbf{y})^k$$

For example, for the typical value of $c = 1$, the inhomogeneous kernel is a weighted sum of the homogeneous polynomial kernels for all powers up to q , that is,

$$(1 + \mathbf{x}^T \mathbf{y})^q = 1 + q \mathbf{x}^T \mathbf{y} + \binom{q}{2} (\mathbf{x}^T \mathbf{y})^2 + \cdots + q (\mathbf{x}^T \mathbf{y})^{q-1} + (\mathbf{x}^T \mathbf{y})^q$$

Example 5.6. Consider the points \mathbf{x}_1 and \mathbf{x}_2 in Figure 5.1.

$$\mathbf{x}_1 = \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} \quad \mathbf{x}_2 = \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix}$$

The homogeneous quadratic kernel is given as

$$K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2)^2 = 50.01^2 = 2501$$

The inhomogeneous quadratic kernel is given as

$$K(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1^T \mathbf{x}_2)^2 = (1 + 50.01)^2 = 51.01^2 = 2602.02$$

For the polynomial kernel it is possible to construct a mapping ϕ from the input to the feature space. Let n_0, n_1, \dots, n_d denote non-negative integers, such that $\sum_{i=0}^d n_i = q$. Further, let $\mathbf{n} = (n_0, n_1, \dots, n_d)$, and let $|\mathbf{n}| = \sum_{i=0}^d n_i = q$. Also, let $\binom{q}{\mathbf{n}}$ denote the multinomial coefficient

$$\binom{q}{\mathbf{n}} = \binom{q}{n_0, n_1, \dots, n_d} = \frac{q!}{n_0! n_1! \dots n_d!}$$

The multinomial expansion of the inhomogeneous kernel is then given as

$$\begin{aligned} K_q(\mathbf{x}, \mathbf{y}) &= (c + \mathbf{x}^T \mathbf{y})^q = \left(c + \sum_{k=1}^d x_k y_k \right)^q = (c + x_1 y_1 + \cdots + x_d y_d)^q \\ &= \sum_{|\mathbf{n}|=q} \binom{q}{\mathbf{n}} c^{n_0} (x_1 y_1)^{n_1} (x_2 y_2)^{n_2} \dots (x_d y_d)^{n_d} \\ &= \sum_{|\mathbf{n}|=q} \binom{q}{\mathbf{n}} c^{n_0} (x_1^{n_1} x_2^{n_2} \dots x_d^{n_d}) (y_1^{n_1} y_2^{n_2} \dots y_d^{n_d}) \\ &= \sum_{|\mathbf{n}|=q} \left(\sqrt{a_{\mathbf{n}}} \prod_{k=1}^d x_k^{n_k} \right) \left(\sqrt{a_{\mathbf{n}}} \prod_{k=1}^d y_k^{n_k} \right) \\ &= \phi(\mathbf{x})^T \phi(\mathbf{y}) \end{aligned}$$

where $a_{\mathbf{n}} = \binom{q}{\mathbf{n}} c^{n_0}$, and the summation is over all $\mathbf{n} = (n_0, n_1, \dots, n_d)$ such that $|\mathbf{n}| = n_0 + n_1 + \cdots + n_d = q$. Using the notation $\mathbf{x}^{\mathbf{n}} = \prod_{k=1}^d x_k^{n_k}$, the mapping $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is given as the vector

$$\phi(\mathbf{x}) = (\dots, a_{\mathbf{n}} \mathbf{x}^{\mathbf{n}}, \dots)^T = \left(\dots, \sqrt{\binom{q}{\mathbf{n}} c^{n_0}} \prod_{k=1}^d x_k^{n_k}, \dots \right)^T$$

where the variable $\mathbf{n} = (n_0, \dots, n_d)$ ranges over all the possible assignments, such that $|\mathbf{n}| = q$. It can be shown that the dimensionality of the feature space is given as

$$m = \binom{d+q}{q}$$

Example 5.7 (Quadratic Polynomial Kernel). Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ and let $c = 1$. The inhomogeneous quadratic polynomial kernel is given as

$$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^2 = (1 + x_1 y_1 + x_2 y_2)^2$$

The set of all assignments $\mathbf{n} = (n_0, n_1, n_2)$, such that $|\mathbf{n}| = q = 2$, and the corresponding terms in the multinomial expansion are shown below.

Assignments $\mathbf{n} = (n_0, n_1, n_2)$	Coefficient $a_{\mathbf{n}} = \binom{q}{\mathbf{n}} c^{n_0}$	Variables $\mathbf{x}^{\mathbf{n}} \mathbf{y}^{\mathbf{n}} = \prod_{k=1}^d (x_i y_i)^{n_i}$
(1, 1, 0)	2	$x_1 y_1$
(1, 0, 1)	2	$x_2 y_2$
(0, 1, 1)	2	$x_1 y_1 x_2 y_2$
(2, 0, 0)	1	1
(0, 2, 0)	1	$(x_1 y_1)^2$
(0, 0, 2)	1	$(x_2 y_2)^2$

Thus, the kernel can be written as

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 y_1 x_2 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 \\ &= \left(1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2\right) \left(1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2, y_1^2, y_2^2\right)^T \\ &= \phi(\mathbf{x})^T \phi(\mathbf{y}) \end{aligned}$$

When the input space is \mathbb{R}^2 , the dimensionality of the feature space is given as

$$m = \binom{d+q}{q} = \binom{2+2}{2} = \binom{4}{2} = 6$$

In this case the inhomogeneous quadratic kernel with $c = 1$ corresponds to the mapping $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$, given as

$$\phi(\mathbf{x}) = \left(1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2\right)^T$$

For example, for $\mathbf{x}_1 = (5.9, 3)^T$ and $\mathbf{x}_2 = (6.9, 3.1)^T$, we have

$$\begin{aligned} \phi(\mathbf{x}_1) &= \left(1, \sqrt{2} \cdot 5.9, \sqrt{2} \cdot 3, \sqrt{2} \cdot 5.9 \cdot 3, 5.9^2, 3^2\right)^T \\ &= (1, 8.34, 4.24, 25.03, 34.81, 9)^T \\ \phi(\mathbf{x}_2) &= \left(1, \sqrt{2} \cdot 6.9, \sqrt{2} \cdot 3.1, \sqrt{2} \cdot 6.9 \cdot 3.1, 6.9^2, 3.1^2\right)^T \\ &= (1, 9.76, 4.38, 30.25, 47.61, 9.61)^T \end{aligned}$$

Thus, the inhomogeneous kernel value is

$$\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) = 1 + 81.40 + 18.57 + 757.16 + 1657.30 + 86.49 = 2601.92$$

On the other hand, when the input space is \mathbb{R}^2 , the homogeneous quadratic kernel corresponds to the mapping $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$, defined as

$$\phi(\mathbf{x}) = \left(\sqrt{2}x_1x_2, x_1^2, x_2^2 \right)^T$$

because only the degree 2 terms are considered. For example, for \mathbf{x}_1 and \mathbf{x}_2 , we have

$$\phi(\mathbf{x}_1) = \left(\sqrt{2} \cdot 5.9 \cdot 3, 5.9^2, 3^2 \right)^T = (25.03, 34.81, 9)^T$$

$$\phi(\mathbf{x}_2) = \left(\sqrt{2} \cdot 6.9 \cdot 3.1, 6.9^2, 3.1^2 \right)^T = (30.25, 47.61, 9.61)^T$$

and thus

$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) = 757.16 + 1657.3 + 86.49 = 2500.95$$

These values essentially match those shown in Example 5.6 up to four significant digits.

Gaussian Kernel

The Gaussian kernel, also called the Gaussian radial basis function (RBF) kernel, is defined as

$$K(\mathbf{x}, \mathbf{y}) = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2} \right\} \quad (5.10)$$

where $\sigma > 0$ is the spread parameter that plays the same role as the standard deviation in a normal density function. Note that $K(\mathbf{x}, \mathbf{x}) = 1$, and further that the kernel value is inversely related to the distance between the two points \mathbf{x} and \mathbf{y} .

Example 5.8. Consider again the points \mathbf{x}_1 and \mathbf{x}_2 in Figure 5.1:

$$\mathbf{x}_1 = \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} \quad \mathbf{x}_2 = \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix}$$

The squared distance between them is given as

$$\|\mathbf{x}_1 - \mathbf{x}_2\|^2 = \|(-1, -0.1)^T\|^2 = 1^2 + 0.1^2 = 1.01$$

With $\sigma = 1$, the Gaussian kernel is

$$K(\mathbf{x}_1, \mathbf{x}_2) = \exp \left\{ -\frac{1.01^2}{2} \right\} = \exp\{-0.51\} = 0.6$$

It is interesting to note that a feature space for the Gaussian kernel has infinite dimensionality. To see this, note that the exponential function can be written as the

infinite expansion

$$\exp\{a\} = \sum_{n=0}^{\infty} \frac{a^n}{n!} = 1 + a + \frac{1}{2!}a^2 + \frac{1}{3!}a^3 + \dots$$

Further, using $\gamma = \frac{1}{2\sigma^2}$, and noting that $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\mathbf{x}^T\mathbf{y}$, we can rewrite the Gaussian kernel as follows:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \exp\left\{-\gamma \|\mathbf{x} - \mathbf{y}\|^2\right\} \\ &= \exp\left\{-\gamma \|\mathbf{x}\|^2\right\} \cdot \exp\left\{-\gamma \|\mathbf{y}\|^2\right\} \cdot \exp\left\{2\gamma \mathbf{x}^T\mathbf{y}\right\} \end{aligned}$$

In particular, the last term is given as the infinite expansion

$$\exp\{2\gamma \mathbf{x}^T\mathbf{y}\} = \sum_{q=0}^{\infty} \frac{(2\gamma)^q}{q!} (\mathbf{x}^T\mathbf{y})^q = 1 + (2\gamma)\mathbf{x}^T\mathbf{y} + \frac{(2\gamma)^2}{2!} (\mathbf{x}^T\mathbf{y})^2 + \dots$$

Using the multinomial expansion of $(\mathbf{x}^T\mathbf{y})^q$, we can write the Gaussian kernel as

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \exp\left\{-\gamma \|\mathbf{x}\|^2\right\} \exp\left\{-\gamma \|\mathbf{y}\|^2\right\} \sum_{q=0}^{\infty} \frac{(2\gamma)^q}{q!} \left(\sum_{|\mathbf{n}|=q} \binom{q}{\mathbf{n}} \prod_{k=1}^d (x_k y_k)^{n_k} \right) \\ &= \sum_{q=0}^{\infty} \sum_{|\mathbf{n}|=q} \left(\sqrt{a_{q,\mathbf{n}}} \exp\left\{-\gamma \|\mathbf{x}\|^2\right\} \prod_{k=1}^d x_k^{n_k} \right) \left(\sqrt{a_{q,\mathbf{n}}} \exp\left\{-\gamma \|\mathbf{y}\|^2\right\} \prod_{k=1}^d y_k^{n_k} \right) \\ &= \phi(\mathbf{x})^T \phi(\mathbf{y}) \end{aligned}$$

where $a_{q,\mathbf{n}} = \frac{(2\gamma)^q}{q!} \binom{q}{\mathbf{n}}$, and $\mathbf{n} = (n_1, n_2, \dots, n_d)$, with $|\mathbf{n}| = n_1 + n_2 + \dots + n_d = q$. The mapping into feature space corresponds to the function $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^\infty$

$$\phi(\mathbf{x}) = \left(\dots, \sqrt{\frac{(2\gamma)^q}{q!} \binom{q}{\mathbf{n}}} \exp\left\{-\gamma \|\mathbf{x}\|^2\right\} \prod_{k=1}^d x_k^{n_k}, \dots \right)^T$$

with the dimensions ranging over all degrees $q = 0, \dots, \infty$, and with the variable $\mathbf{n} = (n_1, \dots, n_d)$ ranging over all possible assignments such that $|\mathbf{n}| = q$ for each value of q . Because ϕ maps the input space into an infinite dimensional feature space, we obviously cannot explicitly transform \mathbf{x} into $\phi(\mathbf{x})$, yet computing the Gaussian kernel $K(\mathbf{x}, \mathbf{y})$ is straightforward.

5.3 BASIC KERNEL OPERATIONS IN FEATURE SPACE

Let us look at some of the basic data analysis tasks that can be performed solely via kernels, without instantiating $\phi(\mathbf{x})$.

Norm of a Point

We can compute the norm of a point $\phi(\mathbf{x})$ in feature space as follows:

$$\|\phi(\mathbf{x})\|^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}) = K(\mathbf{x}, \mathbf{x})$$

which implies that $\|\phi(\mathbf{x})\| = \sqrt{K(\mathbf{x}, \mathbf{x})}$.

Distance between Points

The distance between two points $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ can be computed as

$$\begin{aligned}\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 &= \|\phi(\mathbf{x}_i)\|^2 + \|\phi(\mathbf{x}_j)\|^2 - 2\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &= K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}\quad (5.11)$$

which implies that

$$\delta(\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)) = \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\| = \sqrt{K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)}$$

Rearranging Eq. (5.11), we can see that the kernel value can be considered as a measure of the similarity between two points, as

$$\frac{1}{2} (\|\phi(\mathbf{x}_i)\|^2 + \|\phi(\mathbf{x}_j)\|^2 - \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2) = K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Thus, the more the distance $\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|$ between the two points in feature space, the less the kernel value, that is, the less the similarity.

Example 5.9. Consider the two points \mathbf{x}_1 and \mathbf{x}_2 in Figure 5.1:

$$\mathbf{x}_1 = \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} \quad \mathbf{x}_2 = \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix}$$

Assuming the homogeneous quadratic kernel, the norm of $\phi(\mathbf{x}_1)$ can be computed as

$$\|\phi(\mathbf{x}_1)\|^2 = K(\mathbf{x}_1, \mathbf{x}_1) = (\mathbf{x}_1^T \mathbf{x}_1)^2 = 43.81^2 = 1919.32$$

which implies that the norm of the transformed point is $\|\phi(\mathbf{x}_1)\| = \sqrt{43.81^2} = 43.81$.

The distance between $\phi(\mathbf{x}_1)$ and $\phi(\mathbf{x}_2)$ in feature space is given as

$$\begin{aligned}\delta(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2)) &= \sqrt{K(\mathbf{x}_1, \mathbf{x}_1) + K(\mathbf{x}_2, \mathbf{x}_2) - 2K(\mathbf{x}_1, \mathbf{x}_2)} \\ &= \sqrt{1919.32 + 3274.13 - 2 \cdot 2501} = \sqrt{191.45} = 13.84\end{aligned}$$

Mean in Feature Space

The mean of the points in feature space is given as

$$\mu_\phi = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)$$

Because we do not, in general, have access to $\phi(\mathbf{x}_i)$, we cannot explicitly compute the mean point in feature space.

Nevertheless, we can compute the squared norm of the mean as follows:

$$\begin{aligned}
 \|\boldsymbol{\mu}_\phi\|^2 &= \boldsymbol{\mu}_\phi^T \boldsymbol{\mu}_\phi \\
 &= \left(\frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \right)^T \left(\frac{1}{n} \sum_{j=1}^n \phi(\mathbf{x}_j) \right) \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j)
 \end{aligned} \tag{5.12}$$

The above derivation implies that the squared norm of the mean in feature space is simply the average of the values in the kernel matrix \mathbf{K} .

Example 5.10. Consider the five points from Example 5.3, also shown in Figure 5.1. Example 5.4 showed the norm of the mean for the linear kernel. Let us consider the Gaussian kernel with $\sigma = 1$. The Gaussian kernel matrix is given as

$$\mathbf{K} = \begin{pmatrix} 1.00 & 0.60 & 0.78 & 0.42 & 0.72 \\ 0.60 & 1.00 & 0.94 & 0.07 & 0.44 \\ 0.78 & 0.94 & 1.00 & 0.13 & 0.65 \\ 0.42 & 0.07 & 0.13 & 1.00 & 0.23 \\ 0.72 & 0.44 & 0.65 & 0.23 & 1.00 \end{pmatrix}$$

The squared norm of the mean in feature space is therefore

$$\|\boldsymbol{\mu}_\phi\|^2 = \frac{1}{25} \sum_{i=1}^5 \sum_{j=1}^5 K(\mathbf{x}_i, \mathbf{x}_j) = \frac{14.98}{25} = 0.599$$

which implies that $\|\boldsymbol{\mu}_\phi\| = \sqrt{0.599} = 0.774$.

Total Variance in Feature Space

Let us first derive a formula for the squared distance of a point $\phi(\mathbf{x}_i)$ to the mean $\boldsymbol{\mu}_\phi$ in feature space:

$$\begin{aligned}
 \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi\|^2 &= \|\phi(\mathbf{x}_i)\|^2 - 2\phi(\mathbf{x}_i)^T \boldsymbol{\mu}_\phi + \|\boldsymbol{\mu}_\phi\|^2 \\
 &= K(\mathbf{x}_i, \mathbf{x}_i) - \frac{2}{n} \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n^2} \sum_{a=1}^n \sum_{b=1}^n K(\mathbf{x}_a, \mathbf{x}_b)
 \end{aligned}$$

The total variance [Eq. (1.4)] in feature space is obtained by taking the average squared deviation of points from the mean in feature space:

$$\sigma_\phi^2 = \frac{1}{n} \sum_{i=1}^n \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi\|^2$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n \left(K(\mathbf{x}_i, \mathbf{x}_i) - \frac{2}{n} \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n^2} \sum_{a=1}^n \sum_{b=1}^n K(\mathbf{x}_a, \mathbf{x}_b) \right) \\
&= \frac{1}{n} \sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_i) - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) + \frac{n}{n^3} \sum_{a=1}^n \sum_{b=1}^n K(\mathbf{x}_a, \mathbf{x}_b) \\
&= \frac{1}{n} \sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_i) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \tag{5.13}
\end{aligned}$$

In other words, the total variance in feature space is given as the difference between the average of the diagonal entries and the average of the entire kernel matrix \mathbf{K} . Also notice that by Eq. (5.12) the second term is simply $\|\boldsymbol{\mu}_\phi\|^2$.

Example 5.11. Continuing Example 5.10, the total variance in feature space for the five points, for the Gaussian kernel, is given as

$$\sigma_\phi^2 = \left(\frac{1}{n} \sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_i) \right) - \|\boldsymbol{\mu}_\phi\|^2 = \frac{1}{5} \times 5 - 0.599 = 0.401$$

The distance between $\phi(\mathbf{x}_1)$ and the mean $\boldsymbol{\mu}_\phi$ in feature space is given as

$$\begin{aligned}
\|\phi(\mathbf{x}_1) - \boldsymbol{\mu}_\phi\|^2 &= K(\mathbf{x}_1, \mathbf{x}_1) - \frac{2}{5} \sum_{j=1}^5 K(\mathbf{x}_1, \mathbf{x}_j) + \|\boldsymbol{\mu}_\phi\|^2 \\
&= 1 - \frac{2}{5} (1 + 0.6 + 0.78 + 0.42 + 0.72) + 0.599 \\
&= 1 - 1.410 + 0.599 = 0.189
\end{aligned}$$

Centering in Feature Space

We can center each point in feature space by subtracting the mean from it, as follows:

$$\hat{\phi}(\mathbf{x}_i) = \phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi$$

Because we do not have explicit representation of $\phi(\mathbf{x}_i)$ or $\boldsymbol{\mu}_\phi$, we cannot explicitly center the points. However, we can still compute the *centered kernel matrix*, that is, the kernel matrix over centered points.

The centered kernel matrix is given as

$$\hat{\mathbf{K}} = \left\{ \hat{K}(\mathbf{x}_i, \mathbf{x}_j) \right\}_{i,j=1}^n$$

where each cell corresponds to the kernel between centered points, that is

$$\begin{aligned}
\hat{K}(\mathbf{x}_i, \mathbf{x}_j) &= \hat{\phi}(\mathbf{x}_i)^T \hat{\phi}(\mathbf{x}_j) \\
&= (\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi)^T (\phi(\mathbf{x}_j) - \boldsymbol{\mu}_\phi) \\
&= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) - \phi(\mathbf{x}_i)^T \boldsymbol{\mu}_\phi - \phi(\mathbf{x}_j)^T \boldsymbol{\mu}_\phi + \boldsymbol{\mu}_\phi^T \boldsymbol{\mu}_\phi
\end{aligned}$$

$$\begin{aligned}
&= K(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n} \sum_{k=1}^n \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_k) - \frac{1}{n} \sum_{k=1}^n \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_k) + \|\boldsymbol{\mu}_\phi\|^2 \\
&= K(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_i, \mathbf{x}_k) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_j, \mathbf{x}_k) + \frac{1}{n^2} \sum_{a=1}^n \sum_{b=1}^n K(\mathbf{x}_a, \mathbf{x}_b)
\end{aligned}$$

In other words, we can compute the centered kernel matrix using only the kernel function. Over all the pairs of points, the centered kernel matrix can be written compactly as follows:

$$\begin{aligned}
\hat{\mathbf{K}} &= \mathbf{K} - \frac{1}{n} \mathbf{1}_{n \times n} \mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1}_{n \times n} + \frac{1}{n^2} \mathbf{1}_{n \times n} \mathbf{K} \mathbf{1}_{n \times n} \\
&= \left(\mathbf{I} - \frac{1}{n} \mathbf{1}_{n \times n} \right) \mathbf{K} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}_{n \times n} \right)
\end{aligned} \tag{5.14}$$

where $\mathbf{1}_{n \times n}$ is the $n \times n$ singular matrix, all of whose entries equal 1.

Example 5.12. Consider the first five points from the 2-dimensional Iris dataset shown in Figure 5.1a:

$$\mathbf{x}_1 = \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} \quad \mathbf{x}_2 = \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix} \quad \mathbf{x}_3 = \begin{pmatrix} 6.6 \\ 2.9 \end{pmatrix} \quad \mathbf{x}_4 = \begin{pmatrix} 4.6 \\ 3.2 \end{pmatrix} \quad \mathbf{x}_5 = \begin{pmatrix} 6 \\ 2.2 \end{pmatrix}$$

Consider the linear kernel matrix shown in Figure 5.1b. We can center it by first computing

$$\mathbf{I} - \frac{1}{5} \mathbf{1}_{5 \times 5} = \begin{pmatrix} 0.8 & -0.2 & -0.2 & -0.2 & -0.2 \\ -0.2 & 0.8 & -0.2 & -0.2 & -0.2 \\ -0.2 & -0.2 & 0.8 & -0.2 & -0.2 \\ -0.2 & -0.2 & -0.2 & 0.8 & -0.2 \\ -0.2 & -0.2 & -0.2 & -0.2 & 0.8 \end{pmatrix}$$

The centered kernel matrix [Eq. (5.14)] is given as

$$\begin{aligned}
\hat{\mathbf{K}} &= \left(\mathbf{I} - \frac{1}{5} \mathbf{1}_{5 \times 5} \right) \cdot \begin{pmatrix} 43.81 & 50.01 & 47.64 & 36.74 & 42.00 \\ 50.01 & 57.22 & 54.53 & 41.66 & 48.22 \\ 47.64 & 54.53 & 51.97 & 39.64 & 45.98 \\ 36.74 & 41.66 & 39.64 & 31.40 & 34.64 \\ 42.00 & 48.22 & 45.98 & 34.64 & 40.84 \end{pmatrix} \cdot \left(\mathbf{I} - \frac{1}{5} \mathbf{1}_{5 \times 5} \right) \\
&= \begin{pmatrix} 0.02 & -0.06 & -0.06 & 0.18 & -0.08 \\ -0.06 & 0.86 & 0.54 & -1.19 & -0.15 \\ -0.06 & 0.54 & 0.36 & -0.83 & -0.01 \\ 0.18 & -1.19 & -0.83 & 2.06 & -0.22 \\ -0.08 & -0.15 & -0.01 & -0.22 & 0.46 \end{pmatrix}
\end{aligned}$$

To verify that $\hat{\mathbf{K}}$ is the same as the kernel matrix for the centered points, let us first center the points by subtracting the mean $\boldsymbol{\mu} = (6.0, 2.88)^T$. The centered points

in feature space are given as

$$\mathbf{z}_1 = \begin{pmatrix} -0.1 \\ 0.12 \end{pmatrix} \quad \mathbf{z}_2 = \begin{pmatrix} 0.9 \\ 0.22 \end{pmatrix} \quad \mathbf{z}_3 = \begin{pmatrix} 0.6 \\ 0.02 \end{pmatrix} \quad \mathbf{z}_4 = \begin{pmatrix} -1.4 \\ 0.32 \end{pmatrix} \quad \mathbf{z}_5 = \begin{pmatrix} 0.0 \\ -0.68 \end{pmatrix}$$

For example, the kernel between $\phi(\mathbf{z}_1)$ and $\phi(\mathbf{z}_2)$ is

$$\phi(\mathbf{z}_1)^T \phi(\mathbf{z}_2) = \mathbf{z}_1^T \mathbf{z}_2 = -0.09 + 0.03 = -0.06$$

which matches $\hat{\mathbf{K}}(\mathbf{x}_1, \mathbf{x}_2)$, as expected. The other entries can be verified in a similar manner. Thus, the kernel matrix obtained by centering the data and then computing the kernel is the same as that obtained via Eq. (5.14).

Normalizing in Feature Space

A common form of normalization is to ensure that points in feature space have unit length by replacing $\phi(\mathbf{x}_i)$ with the corresponding unit vector $\phi_n(\mathbf{x}_i) = \frac{\phi(\mathbf{x}_i)}{\|\phi(\mathbf{x}_i)\|}$. The dot product in feature space then corresponds to the cosine of the angle between the two mapped points, because

$$\phi_n(\mathbf{x}_i)^T \phi_n(\mathbf{x}_j) = \frac{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_i)\| \cdot \|\phi(\mathbf{x}_j)\|} = \cos \theta$$

If the mapped points are both centered and normalized, then a dot product corresponds to the correlation between the two points in feature space.

The normalized kernel matrix, \mathbf{K}_n , can be computed using only the kernel function K , as

$$\mathbf{K}_n(\mathbf{x}_i, \mathbf{x}_j) = \frac{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)}{\|\phi(\mathbf{x}_i)\| \cdot \|\phi(\mathbf{x}_j)\|} = \frac{K(\mathbf{x}_i, \mathbf{x}_j)}{\sqrt{K(\mathbf{x}_i, \mathbf{x}_i) \cdot K(\mathbf{x}_j, \mathbf{x}_j)}}$$

\mathbf{K}_n has all diagonal elements as 1.

Let \mathbf{W} denote the diagonal matrix comprising the diagonal elements of \mathbf{K} :

$$\mathbf{W} = \text{diag}(\mathbf{K}) = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & 0 & \cdots & 0 \\ 0 & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

The normalized kernel matrix can then be expressed compactly as

$$\mathbf{K}_n = \mathbf{W}^{-1/2} \cdot \mathbf{K} \cdot \mathbf{W}^{-1/2}$$

where $\mathbf{W}^{-1/2}$ is the diagonal matrix, defined as $\mathbf{W}^{-1/2}(\mathbf{x}_i, \mathbf{x}_i) = \frac{1}{\sqrt{K(\mathbf{x}_i, \mathbf{x}_i)}}$, with all other elements being zero.

Example 5.13. Consider the five points and the linear kernel matrix shown in Figure 5.1. We have

$$\mathbf{W} = \begin{pmatrix} 43.81 & 0 & 0 & 0 & 0 \\ 0 & 57.22 & 0 & 0 & 0 \\ 0 & 0 & 51.97 & 0 & 0 \\ 0 & 0 & 0 & 31.40 & 0 \\ 0 & 0 & 0 & 0 & 40.84 \end{pmatrix}$$

The normalized kernel is given as

$$\mathbf{K}_n = \mathbf{W}^{-1/2} \cdot \mathbf{K} \cdot \mathbf{W}^{-1/2} = \begin{pmatrix} 1.0000 & 0.9988 & 0.9984 & 0.9906 & 0.9929 \\ 0.9988 & 1.0000 & 0.9999 & 0.9828 & 0.9975 \\ 0.9984 & 0.9999 & 1.0000 & 0.9812 & 0.9980 \\ 0.9906 & 0.9828 & 0.9812 & 1.0000 & 0.9673 \\ 0.9929 & 0.9975 & 0.9980 & 0.9673 & 1.0000 \end{pmatrix}$$

The same kernel is obtained if we first normalize the feature vectors to have unit length and then take the dot products. For example, with the linear kernel, the normalized point $\phi_n(\mathbf{x}_1)$ is given as

$$\phi_n(\mathbf{x}_1) = \frac{\phi(\mathbf{x}_1)}{\|\phi(\mathbf{x}_1)\|} = \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|} = \frac{1}{\sqrt{43.81}} \begin{pmatrix} 5.9 \\ 3 \end{pmatrix} = \begin{pmatrix} 0.8914 \\ 0.4532 \end{pmatrix}$$

Likewise, we have $\phi_n(\mathbf{x}_2) = \frac{1}{\sqrt{57.22}} \begin{pmatrix} 6.9 \\ 3.1 \end{pmatrix} = \begin{pmatrix} 0.9122 \\ 0.4098 \end{pmatrix}$. Their dot product is

$$\phi_n(\mathbf{x}_1)^T \phi_n(\mathbf{x}_2) = 0.8914 \cdot 0.9122 + 0.4532 \cdot 0.4098 = 0.9988$$

which matches $\mathbf{K}_n(\mathbf{x}_1, \mathbf{x}_2)$.

If we start with the centered kernel matrix $\hat{\mathbf{K}}$ from Example 5.12, and then normalize it, we obtain the normalized and centered kernel matrix $\hat{\mathbf{K}}_n$:

$$\hat{\mathbf{K}}_n = \begin{pmatrix} 1.00 & -0.44 & -0.61 & 0.80 & -0.77 \\ -0.44 & 1.00 & 0.98 & -0.89 & -0.24 \\ -0.61 & 0.98 & 1.00 & -0.97 & -0.03 \\ 0.80 & -0.89 & -0.97 & 1.00 & -0.22 \\ -0.77 & -0.24 & -0.03 & -0.22 & 1.00 \end{pmatrix}$$

As noted earlier, the kernel value $\hat{\mathbf{K}}_n(\mathbf{x}_i, \mathbf{x}_j)$ denotes the correlation between \mathbf{x}_i and \mathbf{x}_j in feature space, that is, it is cosine of the angle between the centered points $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$.

5.4 KERNELS FOR COMPLEX OBJECTS

We conclude this chapter with some examples of kernels defined for complex data such as strings and graphs. The use of kernels for dimensionality reduction is described in

Section 7.3, for clustering in Section 13.2 and Chapter 16, for discriminant analysis in Section 20.2, and for classification in Sections 21.4 and 21.5.

5.4.1 Spectrum Kernel for Strings

Consider text or sequence data defined over an alphabet Σ . The l -spectrum feature map is the mapping $\phi : \Sigma^* \rightarrow \mathbb{R}^{|\Sigma|^l}$ from the set of substrings over Σ to the $|\Sigma|^l$ -dimensional space representing the number of occurrences of all possible substrings of length l , defined as

$$\phi(\mathbf{x}) = \left(\dots, \#(\alpha), \dots \right)_{\alpha \in \Sigma^l}^T$$

where $\#(\alpha)$ is the number of occurrences of the l -length string α in \mathbf{x} .

The (full) spectrum map is an extension of the l -spectrum map, obtained by considering all lengths from $l = 0$ to $l = \infty$, leading to an infinite dimensional feature map $\phi : \Sigma^* \rightarrow \mathbb{R}^\infty$:

$$\phi(\mathbf{x}) = \left(\dots, \#(\alpha), \dots \right)_{\alpha \in \Sigma^*}^T$$

where $\#(\alpha)$ is the number of occurrences of the string α in \mathbf{x} .

The (l -)spectrum kernel between two strings $\mathbf{x}_i, \mathbf{x}_j$ is simply the dot product between their (l -)spectrum maps:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

A naive computation of the l -spectrum kernel takes $O(|\Sigma|^l)$ time. However, for a given string \mathbf{x} of length n , the vast majority of the l -length strings have an occurrence count of zero, which can be ignored. The l -spectrum map can be effectively computed in $O(n)$ time for a string of length n (assuming $n \gg l$) because there can be at most $n - l + 1$ substrings of length l , and the l -spectrum kernel can thus be computed in $O(n + m)$ time for any two strings of length n and m , respectively.

The feature map for the (full) spectrum kernel is infinite dimensional, but once again, for a given string \mathbf{x} of length n , the vast majority of the strings will have an occurrence count of zero. A straightforward implementation of the spectrum map for a string \mathbf{x} of length n can be computed in $O(n^2)$ time because \mathbf{x} can have at most $\sum_{l=1}^n n - l + 1 = n(n + 1)/2$ distinct nonempty substrings. The spectrum kernel can then be computed in $O(n^2 + m^2)$ time for any two strings of length n and m , respectively. However, a much more efficient computation is enabled via suffix trees (see Chapter 10), with a total time of $O(n + m)$.

Example 5.14. Consider sequences over the DNA alphabet $\Sigma = \{A, C, G, T\}$. Let $\mathbf{x}_1 = ACAGCAGTA$, and let $\mathbf{x}_2 = AGCAAGCGAG$. For $l = 3$, the feature space has dimensionality $|\Sigma|^l = 4^3 = 64$. Nevertheless, we do not have to map the input points into the full feature space; we can compute the reduced 3-spectrum mapping by counting the number of occurrences for only the length 3 substrings that occur in each input sequence, as follows:

$$\phi(\mathbf{x}_1) = (ACA : 1, AGC : 1, AGT : 1, CAG : 2, GCA : 1, GTA : 1)$$

$$\phi(\mathbf{x}_2) = (AAG : 1, AGC : 2, CAA : 1, CGA : 1, GAG : 1, GCA : 1, GCG : 1)$$

where the notation $\alpha : \#(\alpha)$ denotes that substring α has $\#(\alpha)$ occurrences in \mathbf{x}_i . We can then compute the dot product by considering only the common substrings, as follows:

$$K(\mathbf{x}_1, \mathbf{x}_2) = 1 \times 2 + 1 \times 1 = 2 + 1 = 3$$

The first term in the dot product is due to the substring AGC , and the second is due to GCA , which are the only common length 3 substrings between \mathbf{x}_1 and \mathbf{x}_2 .

The full spectrum can be computed by considering the occurrences of all common substrings over all possible lengths. For \mathbf{x}_1 and \mathbf{x}_2 , the common substrings and their occurrence counts are given as

α	A	C	G	AG	CA	AGC	GCA	$AGCA$
$\#(\alpha)$ in \mathbf{x}_1	4	2	2	2	2	1	1	1
$\#(\alpha)$ in \mathbf{x}_2	4	2	4	3	1	2	1	1

Thus, the full spectrum kernel value is given as

$$K(\mathbf{x}_1, \mathbf{x}_2) = 16 + 4 + 8 + 6 + 2 + 2 + 1 + 1 = 40$$

5.4.2 Diffusion Kernels on Graph Nodes

Let \mathbf{S} be some symmetric similarity matrix between nodes of a graph $G = (V, E)$. For instance, \mathbf{S} can be the (weighted) adjacency matrix \mathbf{A} [Eq. (4.1)] or the Laplacian matrix $\mathbf{L} = \mathbf{A} - \mathbf{\Delta}$ (or its negation), where $\mathbf{\Delta}$ is the degree matrix for an undirected graph G , defined as $\Delta(i, i) = d_i$ and $\Delta(i, j) = 0$ for all $i \neq j$, and d_i is the degree of node i .

Consider the similarity between any two nodes obtained by summing the product of the similarities over paths of length 2:

$$S^{(2)}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{a=1}^n S(\mathbf{x}_i, \mathbf{x}_a) S(\mathbf{x}_a, \mathbf{x}_j) = \mathbf{S}_i^T \mathbf{S}_j$$

where

$$\mathbf{S}_i = \left(S(\mathbf{x}_i, \mathbf{x}_1), S(\mathbf{x}_i, \mathbf{x}_2), \dots, S(\mathbf{x}_i, \mathbf{x}_n) \right)^T$$

denotes the (column) vector representing the i -th row of \mathbf{S} (and because \mathbf{S} is symmetric, it also denotes the i th column of \mathbf{S}). Over all pairs of nodes the similarity matrix over paths of length 2, denoted $\mathbf{S}^{(2)}$, is thus given as the square of the base similarity matrix \mathbf{S} :

$$\mathbf{S}^{(2)} = \mathbf{S} \times \mathbf{S} = \mathbf{S}^2$$

In general, if we sum up the product of the base similarities over all l -length paths between two nodes, we obtain the l -length similarity matrix $\mathbf{S}^{(l)}$, which is simply the l th power of \mathbf{S} , that is,

$$\mathbf{S}^{(l)} = \mathbf{S}^l$$

Power Kernels

Even path lengths lead to positive semidefinite kernels, but odd path lengths are not guaranteed to do so, unless the base matrix \mathbf{S} is itself a positive semidefinite matrix. In particular, $\mathbf{K} = \mathbf{S}^2$ is a valid kernel. To see this, assume that the i th row of \mathbf{S} denotes the feature map for \mathbf{x}_i , that is, $\phi(\mathbf{x}_i) = \mathbf{S}_i$. The kernel value between any two points is then a dot product in feature space:

$$K(\mathbf{x}_i, \mathbf{x}_j) = S^{(2)}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{S}_i^T \mathbf{S}_j = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

For a general path length l , let $\mathbf{K} = \mathbf{S}^l$. Consider the eigen-decomposition of \mathbf{S} :

$$\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \sum_{i=1}^n \mathbf{u}_i \lambda_i \mathbf{u}_i^T$$

where \mathbf{U} is the orthogonal matrix of eigenvectors and $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues of \mathbf{S} :

$$\mathbf{U} = \begin{pmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & & | \end{pmatrix} \quad \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{pmatrix}$$

The eigen-decomposition of \mathbf{K} can be obtained as follows:

$$\mathbf{K} = \mathbf{S}^l = (\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T)^l = \mathbf{U} (\mathbf{\Lambda}^l) \mathbf{U}^T$$

where we used the fact that eigenvectors of \mathbf{S} and \mathbf{S}^l are identical, and further that eigenvalues of \mathbf{S}^l are given as $(\lambda_i)^l$ (for all $i = 1, \dots, n$), where λ_i is an eigenvalue of \mathbf{S} . For $\mathbf{K} = \mathbf{S}^l$ to be a positive semidefinite matrix, all its eigenvalues must be non-negative, which is guaranteed for all even path lengths. Because $(\lambda_i)^l$ will be negative if l is odd and λ_i is negative, odd path lengths lead to a positive semidefinite kernel only if \mathbf{S} is positive semidefinite.

Exponential Diffusion Kernel

Instead of fixing the path length *a priori*, we can obtain a new kernel between nodes of a graph by considering paths of all possible lengths, but by damping the contribution of longer paths, which leads to the *exponential diffusion kernel*, defined as

$$\begin{aligned} \mathbf{K} &= \sum_{l=0}^{\infty} \frac{1}{l!} \beta^l \mathbf{S}^l \\ &= \mathbf{I} + \beta \mathbf{S} + \frac{1}{2!} \beta^2 \mathbf{S}^2 + \frac{1}{3!} \beta^3 \mathbf{S}^3 + \cdots \\ &= \exp\{\beta \mathbf{S}\} \end{aligned} \tag{5.15}$$

where β is a damping factor, and $\exp\{\beta \mathbf{S}\}$ is the matrix exponential. The series on the right hand side above converges for all $\beta \geq 0$.

Substituting $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ in Eq. (5.15), and utilizing the fact that $\mathbf{U}\mathbf{U}^T = \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T = \mathbf{I}$, we have

$$\begin{aligned}
 \mathbf{K} &= \mathbf{I} + \beta \mathbf{S} + \frac{1}{2!} \beta^2 \mathbf{S}^2 + \dots \\
 &= \left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T \right) + \left(\sum_{i=1}^n \mathbf{u}_i \beta \lambda_i \mathbf{u}_i^T \right) + \left(\sum_{i=1}^n \mathbf{u}_i \frac{1}{2!} \beta^2 \lambda_i^2 \mathbf{u}_i^T \right) + \dots \\
 &= \sum_{i=1}^n \mathbf{u}_i \left(1 + \beta \lambda_i + \frac{1}{2!} \beta^2 \lambda_i^2 + \dots \right) \mathbf{u}_i^T \\
 &= \sum_{i=1}^n \mathbf{u}_i \exp\{\beta \lambda_i\} \mathbf{u}_i^T \\
 &= \mathbf{U} \begin{pmatrix} \exp\{\beta \lambda_1\} & 0 & \dots & 0 \\ 0 & \exp\{\beta \lambda_2\} & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \exp\{\beta \lambda_n\} \end{pmatrix} \mathbf{U}^T \tag{5.16}
 \end{aligned}$$

Thus, the eigenvectors of \mathbf{K} are the same as those for \mathbf{S} , whereas its eigenvalues are given as $\exp\{\beta \lambda_i\}$, where λ_i is an eigenvalue of \mathbf{S} . Further, \mathbf{K} is symmetric because \mathbf{S} is symmetric, and its eigenvalues are real and non-negative because the exponential of a real number is non-negative. \mathbf{K} is thus a positive semidefinite kernel matrix. The complexity of computing the diffusion kernel is $O(n^3)$ corresponding to the complexity of computing the eigen-decomposition.

Von Neumann Diffusion Kernel

A related kernel based on powers of \mathbf{S} is the *von Neumann diffusion kernel*, defined as

$$\mathbf{K} = \sum_{l=0}^{\infty} \beta^l \mathbf{S}^l \tag{5.17}$$

where $\beta \geq 0$. Expanding Eq. (5.17), we have

$$\begin{aligned}
 \mathbf{K} &= \mathbf{I} + \beta \mathbf{S} + \beta^2 \mathbf{S}^2 + \beta^3 \mathbf{S}^3 + \dots \\
 &= \mathbf{I} + \beta \mathbf{S} (\mathbf{I} + \beta \mathbf{S} + \beta^2 \mathbf{S}^2 + \dots) \\
 &= \mathbf{I} + \beta \mathbf{S} \mathbf{K}
 \end{aligned}$$

Rearranging the terms in the preceding equation, we obtain a closed form expression for the von Neumann kernel:

$$\begin{aligned}
 \mathbf{K} - \beta \mathbf{S} \mathbf{K} &= \mathbf{I} \\
 (\mathbf{I} - \beta \mathbf{S}) \mathbf{K} &= \mathbf{I} \\
 \mathbf{K} &= (\mathbf{I} - \beta \mathbf{S})^{-1} \tag{5.18}
 \end{aligned}$$

Plugging in the eigen-decomposition $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, and rewriting $\mathbf{I} = \mathbf{U}\mathbf{U}^T$, we have

$$\begin{aligned}\mathbf{K} &= \left(\mathbf{U}\mathbf{U}^T - \mathbf{U}(\beta\mathbf{\Lambda})\mathbf{U}^T \right)^{-1} \\ &= \left(\mathbf{U}(\mathbf{I} - \beta\mathbf{\Lambda})\mathbf{U}^T \right)^{-1} \\ &= \mathbf{U}(\mathbf{I} - \beta\mathbf{\Lambda})^{-1}\mathbf{U}^T\end{aligned}$$

where $(\mathbf{I} - \beta\mathbf{\Lambda})^{-1}$ is the diagonal matrix whose i th diagonal entry is $(1 - \beta\lambda_i)^{-1}$. The eigenvectors of \mathbf{K} and \mathbf{S} are identical, but the eigenvalues of \mathbf{K} are given as $1/(1 - \beta\lambda_i)$. For \mathbf{K} to be a positive semidefinite kernel, all its eigenvalues should be non-negative, which in turn implies that

$$\begin{aligned}(1 - \beta\lambda_i)^{-1} &\geq 0 \\ 1 - \beta\lambda_i &\geq 0 \\ \beta &\leq 1/\lambda_i\end{aligned}$$

Further, the inverse matrix $(\mathbf{I} - \beta\mathbf{\Lambda})^{-1}$ exists only if

$$\det(\mathbf{I} - \beta\mathbf{\Lambda}) = \prod_{i=1}^n (1 - \beta\lambda_i) \neq 0$$

which implies that $\beta \neq 1/\lambda_i$ for all i . Thus, for \mathbf{K} to be a valid kernel, we require that $\beta < 1/\lambda_i$ for all $i = 1, \dots, n$. The von Neumann kernel is therefore guaranteed to be positive semidefinite if $|\beta| < 1/\rho(\mathbf{S})$, where $\rho(\mathbf{S}) = \max_i \{|\lambda_i|\}$ is called the *spectral radius* of \mathbf{S} , defined as the largest eigenvalue of \mathbf{S} in absolute value.

Example 5.15. Consider the graph in Figure 5.2. Its adjacency and degree matrices are given as

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix} \quad \mathbf{\Delta} = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{pmatrix}$$

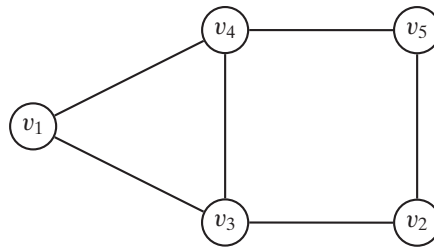


Figure 5.2. Graph diffusion kernel.

The negated Laplacian matrix for the graph is therefore

$$\mathbf{S} = -\mathbf{L} = \mathbf{A} - \mathbf{D} = \begin{pmatrix} -2 & 0 & 1 & 1 & 0 \\ 0 & -2 & 1 & 0 & 1 \\ 1 & 1 & -3 & 1 & 0 \\ 1 & 0 & 1 & -3 & 1 \\ 0 & 1 & 0 & 1 & -2 \end{pmatrix}$$

The eigenvalues of \mathbf{S} are as follows:

$$\lambda_1 = 0 \quad \lambda_2 = -1.38 \quad \lambda_3 = -2.38 \quad \lambda_4 = -3.62 \quad \lambda_5 = -4.62$$

and the eigenvectors of \mathbf{S} are

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \mathbf{u}_4 & \mathbf{u}_5 \\ 0.45 & -0.63 & 0.00 & 0.63 & 0.00 \\ 0.45 & 0.51 & -0.60 & 0.20 & -0.37 \\ 0.45 & -0.20 & -0.37 & -0.51 & 0.60 \\ 0.45 & -0.20 & 0.37 & -0.51 & -0.60 \\ 0.45 & 0.51 & 0.60 & 0.20 & 0.37 \end{pmatrix}$$

Assuming $\beta = 0.2$, the exponential diffusion kernel matrix is given as

$$\begin{aligned} \mathbf{K} = \exp\{0.2\mathbf{S}\} &= \mathbf{U} \begin{pmatrix} \exp\{0.2\lambda_1\} & 0 & \cdots & 0 \\ 0 & \exp\{0.2\lambda_2\} & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & \exp\{0.2\lambda_n\} \end{pmatrix} \mathbf{U}^T \\ &= \begin{pmatrix} 0.70 & 0.01 & 0.14 & 0.14 & 0.01 \\ 0.01 & 0.70 & 0.13 & 0.03 & 0.14 \\ 0.14 & 0.13 & 0.59 & 0.13 & 0.03 \\ 0.14 & 0.03 & 0.13 & 0.59 & 0.13 \\ 0.01 & 0.14 & 0.03 & 0.13 & 0.70 \end{pmatrix} \end{aligned}$$

For the von Neumann diffusion kernel, we have

$$(\mathbf{I} - 0.2\mathbf{A})^{-1} = \begin{pmatrix} 1 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0 & 0.78 & 0.00 & 0.00 & 0.00 \\ 0 & 0.00 & 0.68 & 0.00 & 0.00 \\ 0 & 0.00 & 0.00 & 0.58 & 0.00 \\ 0 & 0.00 & 0.00 & 0.00 & 0.52 \end{pmatrix}$$

For instance, because $\lambda_2 = -1.38$, we have $1 - \beta\lambda_2 = 1 + 0.2 \times 1.38 = 1.28$, and therefore the second diagonal entry is $(1 - \beta\lambda_2)^{-1} = 1/1.28 = 0.78$. The von Neumann kernel is given as

$$\mathbf{K} = \mathbf{U}(\mathbf{I} - 0.2\mathbf{A})^{-1}\mathbf{U}^T = \begin{pmatrix} 0.75 & 0.02 & 0.11 & 0.11 & 0.02 \\ 0.02 & 0.74 & 0.10 & 0.03 & 0.11 \\ 0.11 & 0.10 & 0.66 & 0.10 & 0.03 \\ 0.11 & 0.03 & 0.10 & 0.66 & 0.10 \\ 0.02 & 0.11 & 0.03 & 0.10 & 0.74 \end{pmatrix}$$

5.5 FURTHER READING

Kernel methods have been extensively studied in machine learning and data mining. For an in-depth introduction and more advanced topics see Schölkopf and Smola (2002) and Shawe-Taylor and Cristianini (2004). For applications of kernel methods in bioinformatics see Schölkopf, Tsuda, and Vert (2004).

Schölkopf, B. and Smola, A. J. (2002). *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA: MIT Press.

Schölkopf, B., Tsuda, K., and Vert, J.-P. (2004). *Kernel Methods in Computational Biology*. Cambridge, MA: MIT Press.

Shawe-Taylor, J. and Cristianini, N. (2004). *Kernel Methods for Pattern Analysis*. New York: Cambridge University Press.

5.6 EXERCISES

- Q1.** Prove that the dimensionality of the feature space for the inhomogeneous polynomial kernel of degree q is

$$m = \binom{d+q}{q}$$

- Q2.** Consider the data shown in Table 5.1. Assume the following kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|^2$. Compute the kernel matrix \mathbf{K} .

Table 5.1. Dataset for Q2

i	\mathbf{x}_i
\mathbf{x}_1	(4, 2.9)
\mathbf{x}_2	(2.5, 1)
\mathbf{x}_3	(3.5, 4)
\mathbf{x}_4	(2, 2.1)

- Q3.** Show that eigenvectors of \mathbf{S} and \mathbf{S}^l are identical, and further that eigenvalues of \mathbf{S}^l are given as $(\lambda_i)^l$ (for all $i = 1, \dots, n$), where λ_i is an eigenvalue of \mathbf{S} , and \mathbf{S} is some $n \times n$ symmetric similarity matrix.
- Q4.** The von Neumann diffusion kernel is a valid positive semidefinite kernel if $|\beta| < \frac{1}{\rho(\mathbf{S})}$, where $\rho(\mathbf{S})$ is the spectral radius of \mathbf{S} . Can you derive better bounds for cases when $\beta > 0$ and when $\beta < 0$?
- Q5.** Given the three points $\mathbf{x}_1 = (2.5, 1)^T$, $\mathbf{x}_2 = (3.5, 4)^T$, and $\mathbf{x}_3 = (2, 2.1)^T$.
- (a) Compute the kernel matrix for the Gaussian kernel assuming that $\sigma^2 = 5$.
 - (b) Compute the distance of the point $\phi(\mathbf{x}_1)$ from the mean in feature space.
 - (c) Compute the dominant eigenvector and eigenvalue for the kernel matrix from (a).