

6

Kernel Methods

In Chapters 3 and 4, we considered linear parametric models for regression and classification in which the form of the mapping $y(\mathbf{x}, \mathbf{w})$ from input \mathbf{x} to output y is governed by a vector \mathbf{w} of adaptive parameters. During the learning phase, a set of training data is used either to obtain a point estimate of the parameter vector or to determine a posterior distribution over this vector. The training data is then discarded, and predictions for new inputs are based purely on the learned parameter vector \mathbf{w} . This approach is also used in nonlinear parametric models such as neural networks.

Chapter 5

Section 2.5.1

However, there is a class of pattern recognition techniques, in which the training data points, or a subset of them, are kept and used also during the prediction phase. For instance, the Parzen probability density model comprised a linear combination of ‘kernel’ functions each one centred on one of the training data points. Similarly, in Section 2.5.2 we introduced a simple technique for classification called nearest neighbours, which involved assigning to each new test vector the same label as the

closest example from the training set. These are examples of *memory-based* methods that involve storing the entire training set in order to make predictions for future data points. They typically require a metric to be defined that measures the similarity of any two vectors in input space, and are generally fast to ‘train’ but slow at making predictions for test data points.

Many linear parametric models can be re-cast into an equivalent ‘dual representation’ in which the predictions are also based on linear combinations of a *kernel function* evaluated at the training data points. As we shall see, for models which are based on a fixed nonlinear *feature space* mapping $\phi(\mathbf{x})$, the kernel function is given by the relation

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}'). \quad (6.1)$$

From this definition, we see that the kernel is a symmetric function of its arguments so that $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$. The kernel concept was introduced into the field of pattern recognition by Aizerman *et al.* (1964) in the context of the method of potential functions, so-called because of an analogy with electrostatics. Although neglected for many years, it was re-introduced into machine learning in the context of large-margin classifiers by Boser *et al.* (1992) giving rise to the technique of *support vector machines*. Since then, there has been considerable interest in this topic, both in terms of theory and applications. One of the most significant developments has been the extension of kernels to handle symbolic objects, thereby greatly expanding the range of problems that can be addressed.

The simplest example of a kernel function is obtained by considering the identity mapping for the feature space in (6.1) so that $\phi(\mathbf{x}) = \mathbf{x}$, in which case $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$. We shall refer to this as the *linear kernel*.

The concept of a kernel formulated as an inner product in a feature space allows us to build interesting extensions of many well-known algorithms by making use of the *kernel trick*, also known as *kernel substitution*. The general idea is that, if we have an algorithm formulated in such a way that the input vector \mathbf{x} enters only in the form of scalar products, then we can replace that scalar product with some other choice of kernel. For instance, the technique of kernel substitution can be applied to principal component analysis in order to develop a nonlinear variant of PCA (Schölkopf *et al.*, 1998). Other examples of kernel substitution include nearest-neighbour classifiers and the kernel Fisher discriminant (Mika *et al.*, 1999; Roth and Steinhage, 2000; Baudat and Anouar, 2000).

There are numerous forms of kernel functions in common use, and we shall encounter several examples in this chapter. Many have the property of being a function only of the difference between the arguments, so that $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$, which are known as *stationary kernels* because they are invariant to translations in input space. A further specialization involves *homogeneous kernels*, also known as *radial basis functions*, which depend only on the magnitude of the distance (typically Euclidean) between the arguments so that $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$.

For recent textbooks on kernel methods, see Schölkopf and Smola (2002), Herbrich (2002), and Shawe-Taylor and Cristianini (2004).

Chapter 7

Section 12.3

Section 6.3

6.1. Dual Representations

Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises naturally. This concept will play an important role when we consider support vector machines in the next chapter. Here we consider a linear regression model whose parameters are determined by minimizing a regularized sum-of-squares error function given by

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (6.2)$$

where $\lambda \geq 0$. If we set the gradient of $J(\mathbf{w})$ with respect to \mathbf{w} equal to zero, we see that the solution for \mathbf{w} takes the form of a linear combination of the vectors $\phi(\mathbf{x}_n)$, with coefficients that are functions of \mathbf{w} , of the form

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \quad (6.3)$$

N×M N:样本点的个数 M:基函数的个数
where Φ is the design matrix, whose n^{th} row is given by $\phi(\mathbf{x}_n)^T$. Here the vector $\mathbf{a} = (a_1, \dots, a_N)^T$, and we have defined

$$a_n = -\frac{1}{\lambda} \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}. \quad (6.4)$$

Instead of working with the parameter vector \mathbf{w} , we can now reformulate the least-squares algorithm in terms of the parameter vector \mathbf{a} , giving rise to a *dual representation*. If we substitute $\mathbf{w} = \Phi^T \mathbf{a}$ into $J(\mathbf{w})$, we obtain

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \quad (6.5)$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$. We now define the *Gram matrix* $\mathbf{K} = \Phi \Phi^T$, which is an $N \times N$ symmetric matrix with elements

$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) \quad (6.6)$$

where we have introduced the *kernel function* $k(\mathbf{x}, \mathbf{x}')$ defined by (6.1). In terms of the Gram matrix, the sum-of-squares error function can be written as

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}. \quad (6.7)$$

~~Setting the gradient of $J(\mathbf{a})$ with respect to \mathbf{a} to zero, we obtain the following solution. Using (6.3) to eliminate \mathbf{w} from (6.4) and solving for \mathbf{a} we obtain~~

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}. \quad (6.8)$$

If we substitute this back into the linear regression model, we obtain the following prediction for a new input \mathbf{x}

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad (6.9)$$

where we have defined the vector $\mathbf{k}(\mathbf{x})$ with elements $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$. Thus we see that the dual formulation allows the solution to the least-squares problem to be expressed entirely in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$. This is known as a dual formulation because, by noting that the solution for \mathbf{a} can be expressed as a linear combination of the elements of $\phi(\mathbf{x})$, we recover the original formulation in terms of the parameter vector \mathbf{w} . Note that the prediction at \mathbf{x} is given by a linear combination of the target values from the training set. In fact, we have already obtained this result, using a slightly different notation, in Section 3.3.3.

the original parameter space formulation & dual formulation
对比：前者在M维计算，后者在N维空间中进行计算，后者在N维样本点空间中进行计算。

{ In the dual formulation, we determine the parameter vector \mathbf{a} by inverting an $N \times N$ matrix, whereas in the original parameter space formulation we had to invert an $M \times M$ matrix in order to determine \mathbf{w} . } Because N is typically much larger than M , the dual formulation does not seem to be particularly useful. However, the advantage of the dual formulation, as we shall see, is that it is expressed entirely in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$. We can therefore work directly in terms of kernels and avoid the explicit introduction of the feature vector $\phi(\mathbf{x})$, which allows us implicitly to use feature spaces of high, even infinite, dimensionality. *核方法的好处*

Exercise 6.1

The existence of a dual representation based on the Gram matrix is a property of many linear models, including the perceptron. In Section 6.4, we will develop a duality between probabilistic linear models for regression and the technique of Gaussian processes. Duality will also play an important role when we discuss support vector machines in Chapter 7.

6.2. Constructing Kernels

In order to exploit kernel substitution, we need to be able to construct valid kernel functions. One approach is to choose a feature space mapping $\phi(\mathbf{x})$ and then use this to find the corresponding kernel, as is illustrated in Figure 6.1. Here the kernel function is defined for a one-dimensional input space by

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x') \quad (6.10)$$

where $\phi_i(x)$ are the basis functions.

An alternative approach is to construct kernel functions directly. In this case, we must ensure that the function we choose is a valid kernel, in other words that it corresponds to a scalar product in some (perhaps infinite dimensional) feature space. As a simple example, consider a kernel function given by

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2. \quad (6.11)$$

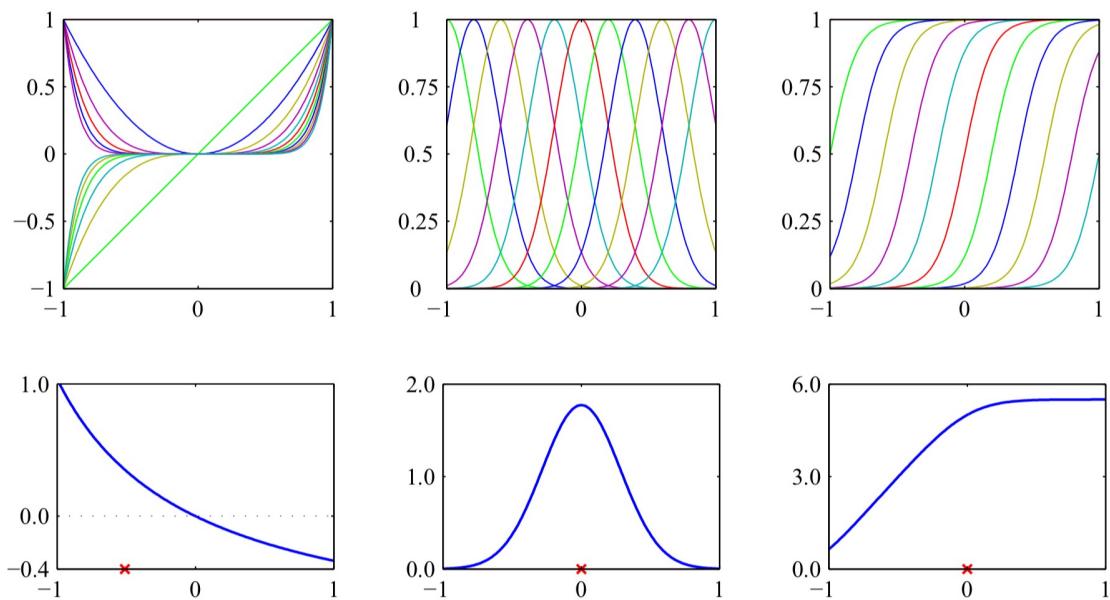


Figure 6.1 Illustration of the construction of kernel functions starting from a corresponding set of basis functions. In each column the lower plot shows the kernel function $k(x, x')$ defined by (6.10) plotted as a function of x , where x' is given by the red cross (\times), while the upper plot shows the corresponding basis functions given by polynomials (left column), ‘Gaussians’ (centre column), and logistic sigmoids (right column).

If we take the particular case of a two-dimensional input space $\mathbf{x} = (x_1, x_2)$ we can expand out the terms and thereby identify the corresponding nonlinear feature mapping

$$\begin{aligned}
 k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\
 &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\
 &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\
 &= \phi(\mathbf{x})^T \phi(\mathbf{z}).
 \end{aligned} \tag{6.12}$$

We see that the feature mapping takes the form $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$ and therefore comprises all possible second order terms, with a specific weighting between them.

More generally, however, we need a simple way to test whether a function constitutes a valid kernel without having to construct the function $\phi(\mathbf{x})$ explicitly. A necessary and sufficient condition for a function $k(\mathbf{x}, \mathbf{x}')$ to be a valid kernel (Shawe-Taylor and Cristianini, 2004) is that the Gram matrix \mathbf{K} , whose elements are given by $k(\mathbf{x}_n, \mathbf{x}_m)$, should be positive semidefinite for all possible choices of the set $\{\mathbf{x}_n\}$. Note that a positive semidefinite matrix is not the same thing as a matrix whose elements are nonnegative.

核函数的必要条件

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\mathbf{x})$ is a function from \mathbf{x} to \mathbb{R}^M , $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M , \mathbf{A} is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

核函数的几个
特点

Equipped with these properties, we can now embark on the construction of more complex kernels appropriate to specific applications. We require that the kernel $k(\mathbf{x}, \mathbf{x}')$ be symmetric and positive semidefinite and that it expresses the appropriate form of similarity between \mathbf{x} and \mathbf{x}' according to the intended application. Here we consider a few common examples of kernel functions. For a more extensive discussion of ‘kernel engineering’, see Shawe-Taylor and Cristianini (2004).

例1： [We saw that the simple polynomial kernel $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$ contains only terms of degree two. If we consider the slightly generalized kernel $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^2$ with $c > 0$, then the corresponding feature mapping $\phi(\mathbf{x})$ contains constant and linear terms as well as terms of order two. Similarly, $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^M$ contains all monomials of order M . For instance, if \mathbf{x} and \mathbf{x}' are two images, then the kernel represents a particular weighted sum of all possible products of M pixels in the first image with M pixels in the second image. This can similarly be generalized to include all terms up to degree M by considering $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$ with $c > 0$. Using the results (6.17) and (6.18) for combining kernels we see that these will all be valid kernel functions.]

例2： [Another commonly used kernel takes the form

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2) \quad (6.23)$$

and is often called a ‘Gaussian’ kernel. Note, however, that in this context it is not interpreted as a probability density, and hence the normalization coefficient is

omitted. We can see that this is a valid kernel by expanding the square

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + (\mathbf{x}')^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}' \quad (6.24)$$

to give

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\mathbf{x}^T \mathbf{x}/2\sigma^2) \exp(\mathbf{x}^T \mathbf{x}'/\sigma^2) \exp(-(\mathbf{x}')^T \mathbf{x}'/2\sigma^2) \quad (6.25)$$

and then making use of (6.14) and (6.16), together with the validity of the linear kernel $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$. Note that the feature vector that corresponds to the Gaussian kernel has infinite dimensionality.

Exercise 6.11

The Gaussian kernel is not restricted to the use of Euclidean distance. If we use kernel substitution in (6.24) to replace $\mathbf{x}^T \mathbf{x}'$ with a nonlinear kernel $\kappa(\mathbf{x}, \mathbf{x}')$, we obtain

$$k(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{1}{2\sigma^2}(\kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}'))\right\}. \quad (6.26)$$

解3： [An important contribution to arise from the kernel viewpoint has been the extension to inputs that are symbolic, rather than simply vectors of real numbers. Kernel functions can be defined over objects as diverse as graphs, sets, strings, and text documents. Consider, for instance, a fixed set and define a nonvectorial space consisting of all possible subsets of this set. If A_1 and A_2 are two such subsets then one simple choice of kernel would be

$$k(A_1, A_2) = 2^{|A_1 \cap A_2|} \quad (6.27)$$

where $A_1 \cap A_2$ denotes the intersection of sets A_1 and A_2 , and $|A|$ denotes the number of elements in A . This is a valid kernel function because it can be shown to correspond to an inner product in a feature space.]

Exercise 6.12

**解4：基于概率模型
构建核函数**

方案一

核①

[One powerful approach to the construction of kernels starts from a probabilistic generative model (Haussler, 1999), which allows us to apply generative models in a discriminative setting. Generative models can deal naturally with missing data and in the case of hidden Markov models can handle sequences of varying length. By contrast, discriminative models generally give better performance on discriminative tasks than generative models. It is therefore of some interest to combine these two approaches (Lasserre *et al.*, 2006). One way to combine them is to use a generative model to define a kernel, and then use this kernel in a discriminative approach.

Given a generative model $p(\mathbf{x})$ we can define a kernel by

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}'). \quad (6.28)$$

This is clearly a valid kernel function because we can interpret it as an inner product in the one-dimensional feature space defined by the mapping $p(\mathbf{x})$. It says that two inputs \mathbf{x} and \mathbf{x}' are similar if they both have high probabilities. We can use (6.13) and (6.17) to extend this class of kernels by considering sums over products of different probability distributions, with positive weighting coefficients $p(i)$, of the form

$$k(\mathbf{x}, \mathbf{x}') = \sum_i p(\mathbf{x}|i)p(\mathbf{x}'|i)p(i). \quad (6.29)$$

核① \Rightarrow 核② \Rightarrow 核③

X, X' 相互独立, 且 $K(x, x') = \sum_i p(x, x' | z) p(z)$, $p(z)$ 是非一阶的模型, 可视为混合模型
 This is equivalent, up to an overall multiplicative constant, to a mixture distribution in which the components factorize, with the index i playing the role of a ‘latent’ variable. Two inputs x and x' will give a large value for the kernel function, and hence appear similar, if they have significant probability under a range of different components. Taking the limit of an infinite sum, we can also consider kernels of the form

$$k(x, x') = \int p(x|z)p(x'|z)p(z) dz \quad (6.30)$$

where z is a continuous latent variable.

Section 9.2

核③

Now suppose that our data consists of ordered sequences of length L so that an observation is given by $\mathbf{X} = \{x_1, \dots, x_L\}$. A popular generative model for sequences is the hidden Markov model, which expresses the distribution $p(\mathbf{X})$ as a marginalization over a corresponding sequence of hidden states $\mathbf{Z} = \{z_1, \dots, z_L\}$. We can use this approach to define a kernel function measuring the similarity of two sequences \mathbf{X} and \mathbf{X}' by extending the mixture representation (6.29) to give

$$k(\mathbf{X}, \mathbf{X}') = \sum_{\mathbf{Z}} p(\mathbf{X}|\mathbf{Z})p(\mathbf{X}'|\mathbf{Z})p(\mathbf{Z}) \quad (6.31)$$

so that both observed sequences are generated by the same hidden sequence \mathbf{Z} . This model can easily be extended to allow sequences of differing length to be compared.

// An alternative technique for using generative models to define kernel functions is known as the *Fisher kernel* (Jaakkola and Haussler, 1999). Consider a parametric generative model $p(x|\theta)$ where θ denotes the vector of parameters. The goal is to find a kernel that measures the similarity of two input vectors x and x' induced by the generative model. Jaakkola and Haussler (1999) consider the gradient with respect to θ , which defines a vector in a ‘feature’ space having the same dimensionality as θ . In particular, they consider the *Fisher score*

$$\mathbf{g}(\theta, x) = \nabla_{\theta} \ln p(x|\theta) \quad (6.32)$$

from which the *Fisher kernel* is defined by

$$k(x, x') = \mathbf{g}(\theta, x)^T \mathbf{F}^{-1} \mathbf{g}(\theta, x'). \quad (6.33)$$

Here \mathbf{F} is the *Fisher information matrix*, given by

$$\mathbf{F} = \mathbb{E}_x [\mathbf{g}(\theta, x)\mathbf{g}(\theta, x)^T] \quad (6.34)$$

where the expectation is with respect to x under the distribution $p(x|\theta)$. This can be motivated from the perspective of *information geometry* (Amari, 1998), which considers the differential geometry of the space of model parameters. Here we simply note that the presence of the Fisher information matrix causes this kernel to be invariant under a nonlinear re-parameterization of the density model $\theta \rightarrow \psi(\theta)$.

In practice, it is often infeasible to evaluate the Fisher information matrix. One approach is simply to replace the expectation in the definition of the Fisher information with the sample average, giving

$$\mathbf{F} \simeq \frac{1}{N} \sum_{n=1}^N \mathbf{g}(\theta, x_n)\mathbf{g}(\theta, x_n)^T. \quad (6.35)$$

Section 13.2

核②的一个例子

上面介绍了使用概率模型构造核函数的第一种方案, 包括核①, ②, ③三种形式;
 下面介绍第二种方案, 即 Fisher kernel

Exercise 6.13

注意, 这里要求 $\psi(\cdot)$ invertible and differentiable

Section 12.1.3

This is the covariance matrix of the Fisher scores, and so the Fisher kernel corresponds to a whitening of these scores. More simply, we can just omit the Fisher information matrix altogether and use the noninvariant kernel

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^T \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}'). \quad (6.36)$$

An application of Fisher kernels to document retrieval is given by Hofmann (2000).]

例子： [A final example of a kernel function is the sigmoidal kernel given by

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a \mathbf{x}^T \mathbf{x}' + b) \quad (6.37)$$

whose Gram matrix in general is not positive semidefinite. This form of kernel has, however, been used in practice (Vapnik, 1995), possibly because it gives kernel expansions such as the support vector machine a superficial resemblance to neural network models. As we shall see, in the limit of an infinite number of basis functions, a Bayesian neural network with an appropriate prior reduces to a Gaussian process, thereby providing a deeper link between neural networks and kernel methods]

Section 6.4.7

6.3. Radial Basis Function Networks

In Chapter 3, we discussed regression models based on linear combinations of fixed basis functions, although we did not discuss in detail what form those basis functions might take. One choice that has been widely used is that of *radial basis functions*, which have the property that each basis function depends only on the radial distance (typically Euclidean) from a centre μ_j , so that $\phi_j(\mathbf{x}) = h(\|\mathbf{x} - \mu_j\|)$.

基向量函数的几个应用
场景或源起

[Historically, radial basis functions were introduced for the purpose of exact function interpolation (Powell, 1987). Given a set of input vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ along with corresponding target values $\{t_1, \dots, t_N\}$, the goal is to find a smooth function $f(\mathbf{x})$ that fits every target value exactly, so that $f(\mathbf{x}_n) = t_n$ for $n = 1, \dots, N$. This is achieved by expressing $f(\mathbf{x})$ as a linear combination of radial basis functions, one centred on every data point

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|). \quad (6.38)$$

The values of the coefficients $\{w_n\}$ are found by least squares, and because there are the same number of coefficients as there are constraints, the result is a function that fits every target value exactly. In pattern recognition applications, however, the target values are generally noisy, and exact interpolation is undesirable because this corresponds to an over-fitted solution.

// Expansions in radial basis functions also arise from regularization theory (Poggio and Girosi, 1990; Bishop, 1995a). For a sum-of-squares error function with a regularizer defined in terms of a differential operator, the optimal solution is given by an expansion in the *Green's functions* of the operator (which are analogous to the eigenvectors of a discrete matrix), again with one basis function centred on each data

point. If the differential operator is isotropic then the Green's functions depend only on the radial distance from the corresponding data point. Due to the presence of the regularizer, the solution no longer interpolates the training data exactly.

// Another motivation for radial basis functions comes from a consideration of the interpolation problem when the input (rather than the target) variables are noisy (Webb, 1994; Bishop, 1995a). If the noise on the input variable \mathbf{x} is described by a variable ξ having a distribution $\nu(\xi)$, then the sum-of-squares error function becomes

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \xi) - t_n\}^2 \nu(\xi) d\xi. \quad (6.39)$$

Appendix D

Exercise 6.17 泛函优化问题 变分法求解

Using the calculus of variations, we can optimize with respect to the function $y(\mathbf{x})$ to give

$$y(\mathbf{x}) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n) \quad (6.40)$$

where the basis functions are given by

$$h(\mathbf{x} - \mathbf{x}_n) = \frac{\nu(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \nu(\mathbf{x} - \mathbf{x}_n)}. \quad (6.41)$$

We see that there is one basis function centred on every data point. This is known as the Nadaraya-Watson model and will be derived again from a different perspective in Section 6.3.1. If the noise distribution $\nu(\xi)$ is isotropic, so that it is a function only of $\|\xi\|$, then the basis functions will be radial.

Note that the basis functions (6.41) are normalized, so that $\sum_n h(\mathbf{x} - \mathbf{x}_n) = 1$ for any value of \mathbf{x} . The effect of such normalization is shown in Figure 6.2. Normalization is sometimes used in practice as it avoids having regions of input space where all of the basis functions take small values, which would necessarily lead to predictions in such regions that are either small or controlled purely by the bias parameter.

// Another situation in which expansions in normalized radial basis functions arise is in the application of kernel density estimation to the problem of regression, as we shall discuss in Section 6.3.1.]

[Because there is one basis function associated with every data point, the corresponding model can be computationally costly to evaluate when making predictions for new data points. Models have therefore been proposed (Broomhead and Lowe, 1988; Moody and Darken, 1989; Poggio and Girosi, 1990), which retain the expansion in radial basis functions but where the number M of basis functions is smaller than the number N of data points. Typically, the number of basis functions, and the locations μ_i of their centres, are determined based on the input data $\{\mathbf{x}_n\}$ alone. The basis functions are then kept fixed and the coefficients $\{w_i\}$ are determined by least squares by solving the usual set of linear equations, as discussed in Section 3.1.1.

径向基函数 中心点的选择

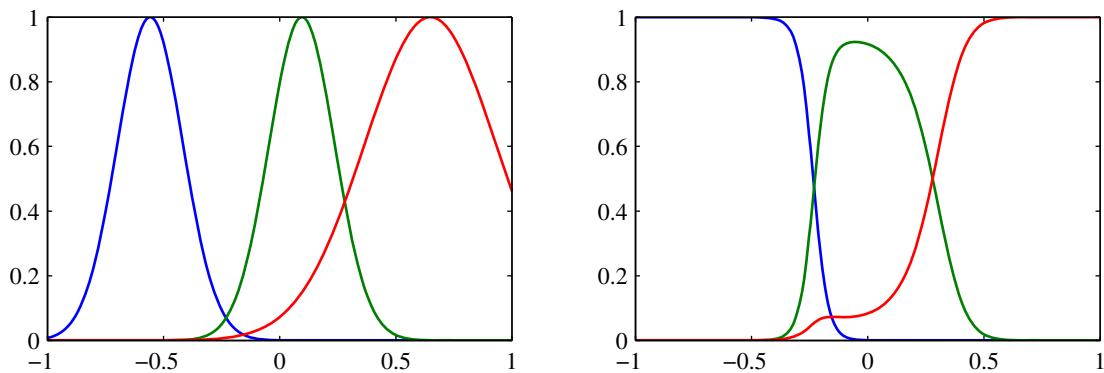


Figure 6.2 Plot of a set of Gaussian basis functions on the left, together with the corresponding normalized basis functions on the right.

三种中心点选取方法 One of the simplest ways of choosing basis function centres is to use a randomly chosen subset of the data points. A more systematic approach is called *orthogonal least squares* (Chen *et al.*, 1991). This is a sequential selection process in which at each step the next data point to be chosen as a basis function centre corresponds to the one that gives the greatest reduction in the sum-of-squares error. Values for the expansion coefficients are determined as part of the algorithm. Clustering algorithms such as *K*-means have also been used, which give a set of basis function centres that no longer coincide with training data points.

Section 9.1

6.3.1 Nadaraya-Watson model

In Section 3.3.3, we saw that the prediction of a linear regression model for a new input \mathbf{x} takes the form of a linear combination of the training set target values with coefficients given by the ‘equivalent kernel’ (3.62) where the equivalent kernel satisfies the summation constraint (3.64).

核回归模型的另一种推导方法
We can motivate the kernel regression model (3.61) from a different perspective, starting with kernel density estimation. Suppose we have a training set $\{\mathbf{x}_n, t_n\}$ and we use a Parzen density estimator to model the joint distribution $p(\mathbf{x}, t)$, so that

定义核密度模型，注意，
这里核密度模型的定义
$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n) \quad (6.42)$$

直接基于样本集 (\mathbf{x}_n, t_n) 而不是参数 f
where $f(\mathbf{x}, t)$ is the component density function, and there is one such component centred on each data point. We now find an expression for the regression function $y(\mathbf{x})$, corresponding to the conditional average of the target variable conditioned on

Section 2.5.1

基于样本空间建模
基于参数空间建模

the input variable, which is given by

根据高斯模型式(6.42)
定义回归方程，
该回归方程实际上
对应最小均方误差
的结果。

$$\begin{aligned}
 y(\mathbf{x}) &= \mathbb{E}[t|\mathbf{x}] = \int_{-\infty}^{\infty} tp(t|\mathbf{x}) dt \\
 &= \frac{\int tp(\mathbf{x}, t) dt}{\int p(\mathbf{x}, t) dt} \\
 &= \frac{\sum_n \int tf(\mathbf{x} - \mathbf{x}_n, t - t_n) dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt}. \tag{6.43}
 \end{aligned}$$

// We now assume for simplicity that the component density functions have zero mean so that

引入假设，恒等变形，最终
得到回归方程基于模型

$$\int_{-\infty}^{\infty} f(\mathbf{x}, t) t dt = 0 \tag{6.44}$$

for all values of \mathbf{x} . Using a simple change of variable, we then obtain

数的表示形式，即式(6.45)

$$\begin{aligned}
 y(\mathbf{x}) &= \frac{\sum_n g(\mathbf{x} - \mathbf{x}_n) t_n}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \\
 &= \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n \tag{6.45}
 \end{aligned}$$

where $n, m = 1, \dots, N$ and the kernel function $k(\mathbf{x}, \mathbf{x}_n)$ is given by

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \tag{6.46}$$

and we have defined

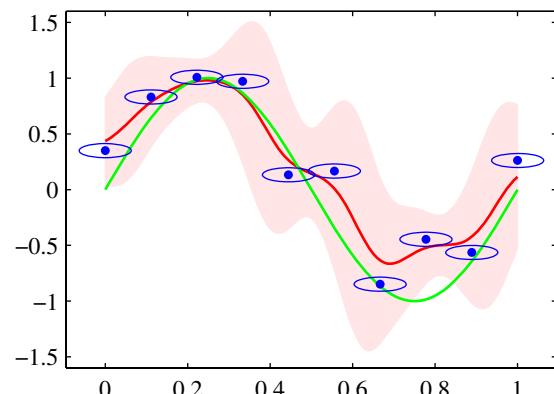
$$g(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}, t) dt. \tag{6.47}$$

The result (6.45) is known as the Nadaraya-Watson model, or *kernel regression* (Nadaraya, 1964; Watson, 1964). For a localized kernel function, it has the property of giving more weight to the data points \mathbf{x}_n that are close to \mathbf{x} . Note that the kernel (6.46) satisfies the summation constraint

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1. \quad \boxed{}$$

Figure 6.3 Illustration of the Nadaraya-Watson kernel regression model using isotropic Gaussian kernels, for the sinusoidal data set. The original sine function is shown by the green curve, the data points are shown in blue, and each is the centre of an isotropic Gaussian kernel. The resulting regression function, given by the conditional mean, is shown by the red line, along with the two-standard-deviation region for the conditional distribution $p(t|x)$ shown by the red shading. The blue ellipse around each data point shows one standard deviation contour for the corresponding kernel. These appear noncircular due to the different scales on the horizontal and vertical axes.

各向同性高斯分布等高线是圆环



// In fact, this model defines not only a conditional expectation but also a full conditional distribution given by

核函数模型(6.42)

对应的条件分布 $p(t|x)$

$$\text{计算公式} \quad p(t|x) = \frac{p(t, x)}{\int p(t, x) dt} = \frac{\sum_n f(\mathbf{x} - \mathbf{x}_n, t - t_n)}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt} \quad (6.48)$$

from which other expectations can be evaluated.]

举个例子 [As an illustration we consider the case of a single input variable x in which $f(x, t)$ is given by a zero-mean isotropic Gaussian over the variable $\mathbf{z} = (x, t)$ with variance σ^2 . The corresponding conditional distribution (6.48) is given by a Gaussian mixture, and is shown, together with the conditional mean, for the sinusoidal synthetic data set in Figure 6.3.

Exercise 6.18

对Figure 6.3的例子
进一步说明

若对 $p(t, x)$ 使用高斯混合模型，则它不再能引入核函数，但它们有自己的优点。实际上GMM是基于参数空间建模而不是像这里基于样本空间建模。

6.4. Gaussian Processes

前面从对偶表示的角度引出了核函数，下面我们将从看到如何从贝叶斯的角度自然地引出核函数。

tic discriminative models, leading to the framework of Gaussian processes. We shall thereby see how kernels arise naturally in a Bayesian setting.

基于参数模型的贝叶斯
回归与高斯过程回归
的不同思路

In Chapter 3, we considered linear regression models of the form $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$ in which \mathbf{w} is a vector of parameters and $\phi(\mathbf{x})$ is a vector of fixed nonlinear basis functions that depend on the input vector \mathbf{x} . We showed that a prior distribution over \mathbf{w} induced a corresponding prior distribution over functions $y(\mathbf{x}, \mathbf{w})$. Given a training data set, we then evaluated the posterior distribution over \mathbf{w} and thereby obtained the corresponding posterior distribution over regression functions, which in turn (with the addition of noise) implies a predictive distribution $p(t|\mathbf{x})$ for new input vectors \mathbf{x} .

参数空间 → 函数空间

In the Gaussian process viewpoint, we dispense with the parametric model and instead define a prior probability distribution over functions directly. At first sight, it might seem difficult to work with a distribution over the uncountably infinite space of functions. However, as we shall see, for a finite training set we only need to consider the values of the function at the discrete set of input values \mathbf{x}_n corresponding to the training set and test set data points, and so in practice we can work in a finite space.

Models equivalent to Gaussian processes have been widely studied in many different fields. For instance, in the geostatistics literature Gaussian process regression is known as *kriging* (Cressie, 1993). Similarly, ARMA (autoregressive moving average) models, Kalman filters, and radial basis function networks can all be viewed as forms of Gaussian process models. Reviews of Gaussian processes from a machine learning perspective can be found in MacKay (1998), Williams (1999), and MacKay (2003), and a comparison of Gaussian process models with alternative approaches is given in Rasmussen (1996). See also Rasmussen and Williams (2006) for a recent textbook on Gaussian processes.

6.4.1 Linear regression revisited

In order to motivate the Gaussian process viewpoint, let us return to the linear regression example and re-derive the predictive distribution by working in terms of distributions over functions $y(\mathbf{x}, \mathbf{w})$. This will provide a specific example of a Gaussian process.

高斯过程的一个例子

Consider a model defined in terms of a linear combination of M fixed basis functions given by the elements of the vector $\phi(\mathbf{x})$ so that

$$\text{注意, 式(6.49)是没有拟合能力的} \quad y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (6.49)$$

where \mathbf{x} is the input vector and \mathbf{w} is the M -dimensional weight vector. Now consider a prior distribution over \mathbf{w} given by an isotropic Gaussian of the form

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (6.50)$$

$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$
 $p(y_i | \mathbf{w}, \mathbf{x}_i) = \begin{cases} 1, & y_i = \mathbf{w}^T \phi(\mathbf{x}_i) \\ 0, & y_i \neq \mathbf{w}^T \phi(\mathbf{x}_i) \end{cases}$

给定 \mathbf{x} , 模型假设 y_i 与 \mathbf{w} 为确定性关系。对于训练样本 $\{(\mathbf{x}_i, y_i)\}$, 似然

$$\sum_{i=1}^n \log p(y_i | \mathbf{w}, \mathbf{x}_i)$$

中每个被加项要满足 (\mathbf{w} 取值使 $y_i = \mathbf{w}^T \phi(\mathbf{x}_i)$). 要

governed by the hyperparameter α , which represents the precision (inverse variance) of the distribution. For any given value of \mathbf{w} , the definition (6.49) defines a particular function of \mathbf{x} . The probability distribution over \mathbf{w} defined by (6.50) therefore induces a probability distribution over functions $y(\mathbf{x})$. In practice, we wish to evaluate this function at specific values of \mathbf{x} , for example at the training data points

均为 $-\infty$ (\mathbf{w} 取值不满足 $y_i = \mathbf{w}^T \phi(\mathbf{x}_i)$)

根本无法最大化, 可见式(6.49)是
没有拟合能力的。

$\mathbf{x}_1, \dots, \mathbf{x}_N$. We are therefore interested in the joint distribution of the function values $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$, which we denote by the vector \mathbf{y} with elements $y_n = y(\mathbf{x}_n)$ for $n = 1, \dots, N$. From (6.49), this vector is given by

$$\mathbf{y} = \Phi \mathbf{w} \quad (6.51)$$

where Φ is the design matrix with elements $\Phi_{nk} = \phi_k(\mathbf{x}_n)$. We can find the probability distribution of \mathbf{y} as follows. First of all we note that \mathbf{y} is a linear combination of Gaussian distributed variables given by the elements of \mathbf{w} and hence is itself Gaussian. We therefore need only to find its mean and covariance, which are given from (6.50) by

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0} \quad (6.52)$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K} \quad (6.53)$$

where \mathbf{K} is the Gram matrix with elements

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \quad (6.54)$$

and $k(\mathbf{x}, \mathbf{x}')$ is the kernel function.]

This model provides us with a particular example of a Gaussian process. In general, a Gaussian process is defined as a probability distribution over functions $y(\mathbf{x})$ such that the set of values of $y(\mathbf{x})$ evaluated at an arbitrary set of points $\mathbf{x}_1, \dots, \mathbf{x}_N$ jointly have a Gaussian distribution. {In cases where the input vector \mathbf{x} is two dimensional, this may also be known as a *Gaussian random field*.} More generally, a stochastic process $y(\mathbf{x})$ is specified by giving the joint probability distribution for any finite set of values $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$ in a consistent manner.

A key point about Gaussian stochastic processes is that the joint distribution over N variables y_1, \dots, y_N is specified completely by the second-order statistics, namely the mean and the covariance. {In most applications, we will not have any prior knowledge about the mean of $y(\mathbf{x})$ and so by symmetry we take it to be zero. This is equivalent to choosing the mean of the prior over weight values $p(\mathbf{w}|\alpha)$ to be zero in the basis function viewpoint.} The specification of the Gaussian process is then completed by giving the covariance of $y(\mathbf{x})$ evaluated at any two values of \mathbf{x} , which is given by the kernel function

$$\mathbb{E}[y(\mathbf{x}_n)y(\mathbf{x}_m)] = k(\mathbf{x}_n, \mathbf{x}_m). \quad (6.55)$$

For the specific case of a Gaussian process defined by the linear regression model (6.49) with a weight prior (6.50), the kernel function is given by (6.54).

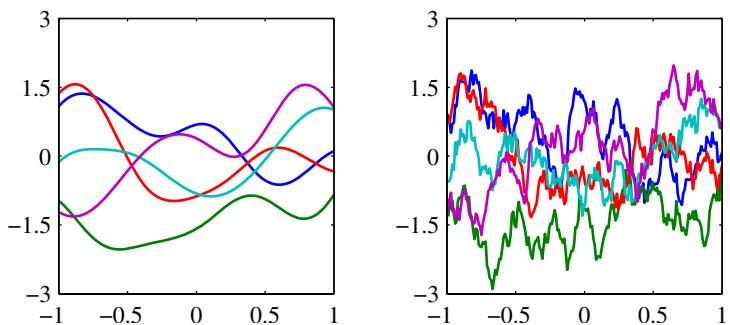
We can also define the kernel function directly, rather than indirectly through a choice of basis function. Figure 6.4 shows samples of functions drawn from Gaussian processes for two different choices of kernel function. The first of these is a ‘Gaussian’ kernel of the form (6.23), and the second is the exponential kernel given by

$$k(x, x') = \exp(-\theta |x - x'|) \quad (6.56)$$

which corresponds to the Ornstein-Uhlenbeck process originally introduced by Uhlenbeck and Ornstein (1930) to describe Brownian motion.

随机过程便和核方法联系了起来，以后我们其位过程的协差矩阵都基于核函数，便可以不用考虑具体基函数向量的选取而直接选择核函数。上述结果是基于式(6.5)这样一个高斯过程模型得到的。

Figure 6.4 Samples from Gaussian processes for a ‘Gaussian’ kernel (left) and an exponential kernel (right).



6.4.2 Gaussian processes for regression

In order to apply Gaussian process models to the problem of regression, we need to take account of the noise on the observed target values, which are given by

$$t_n = y_n + \epsilon_n \quad (6.57)$$

where $y_n = y(\mathbf{x}_n)$, and ϵ_n is a random noise variable whose value is chosen independently for each observation n . Here we shall consider noise processes that have a Gaussian distribution, so that

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}) \quad (6.58)$$

$\epsilon = \{\epsilon_1, \dots, \epsilon_N\}$ 也是一个高斯过程

where β is a hyperparameter representing the precision of the noise. Because the noise is independent for each data point, the joint distribution of the target values $\mathbf{t} = (t_1, \dots, t_N)^T$ conditioned on the values of $\mathbf{y} = (y_1, \dots, y_N)^T$ is given by an isotropic Gaussian of the form

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1} \mathbf{I}_N) \quad (6.59)$$

where \mathbf{I}_N denotes the $N \times N$ unit matrix. From the definition of a Gaussian process, the marginal distribution $p(\mathbf{y})$ is given by a Gaussian whose mean is zero and whose covariance is defined by a Gram matrix \mathbf{K} so that

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}). \quad (6.60)$$

The kernel function that determines \mathbf{K} is typically chosen to express the property that, for points \mathbf{x}_n and \mathbf{x}_m that are similar, the corresponding values $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ will be more strongly correlated than for dissimilar points. Here the notion of similarity will depend on the application.

In order to find the marginal distribution $p(\mathbf{t})$, conditioned on the input values $\mathbf{x}_1, \dots, \mathbf{x}_N$, we need to integrate over \mathbf{y} . This can be done by making use of the results from Section 2.3.3 for the linear-Gaussian model. Using (2.115), we see that the marginal distribution of \mathbf{t} is given by

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (6.61)$$

where the covariance matrix \mathbf{C} has elements

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}. \quad (6.62)$$

This result reflects the fact that the two Gaussian sources of randomness, namely that associated with $y(\mathbf{x})$ and that associated with ϵ , are independent and so their covariances simply add.

One widely used kernel function for Gaussian process regression is given by the exponential of a quadratic form, with the addition of constant and linear terms to give

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left\{ -\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right\} + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m. \quad (6.63)$$

Note that the term involving θ_3 corresponds to a parametric model that is a linear function of the input variables. Samples from this prior are plotted for various values of the parameters $\theta_0, \dots, \theta_3$ in Figure 6.5, and Figure 6.6 shows a set of points sampled from the joint distribution (6.60) along with the corresponding values defined by (6.61).

前面得到了高斯过程的
概率分布，下面将其用于
预测，即计算 $p(t_{N+1} | \mathbf{x})$ // So far, we have used the Gaussian process viewpoint to build a model of the joint distribution over sets of data points. Our goal in regression, however, is to make predictions of the target variables for new inputs, given a set of training data. Let us suppose that $\mathbf{t}_N = (t_1, \dots, t_N)^T$, corresponding to input values $\mathbf{x}_1, \dots, \mathbf{x}_N$, comprise the observed training set, and our goal is to predict the target variable t_{N+1} for a new input vector \mathbf{x}_{N+1} . This requires that we evaluate the predictive distribution $p(t_{N+1} | \mathbf{t}_N)$. Note that this distribution is conditioned also on the variables $\mathbf{x}_1, \dots, \mathbf{x}_N$ and \mathbf{x}_{N+1} . However, to keep the notation simple we will not show these conditioning variables explicitly.

To find the conditional distribution $p(t_{N+1} | \mathbf{t})$, we begin by writing down the joint distribution $p(\mathbf{t}_{N+1})$, where \mathbf{t}_{N+1} denotes the vector $(t_1, \dots, t_N, t_{N+1})^T$. We then apply the results from Section 2.3.1 to obtain the required conditional distribution, as illustrated in Figure 6.7.

From (6.61), the joint distribution over t_1, \dots, t_{N+1} will be given by

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1}) \quad (6.64)$$

where \mathbf{C}_{N+1} is an $(N + 1) \times (N + 1)$ covariance matrix with elements given by (6.62). Because this joint distribution is Gaussian, we can apply the results from Section 2.3.1 to find the conditional Gaussian distribution. To do this, we partition the covariance matrix as follows

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix} \quad (6.65)$$

where \mathbf{C}_N is the $N \times N$ covariance matrix with elements given by (6.62) for $n, m = 1, \dots, N$, the vector \mathbf{k} has elements $k(\mathbf{x}_n, \mathbf{x}_{N+1})$ for $n = 1, \dots, N$, and the scalar

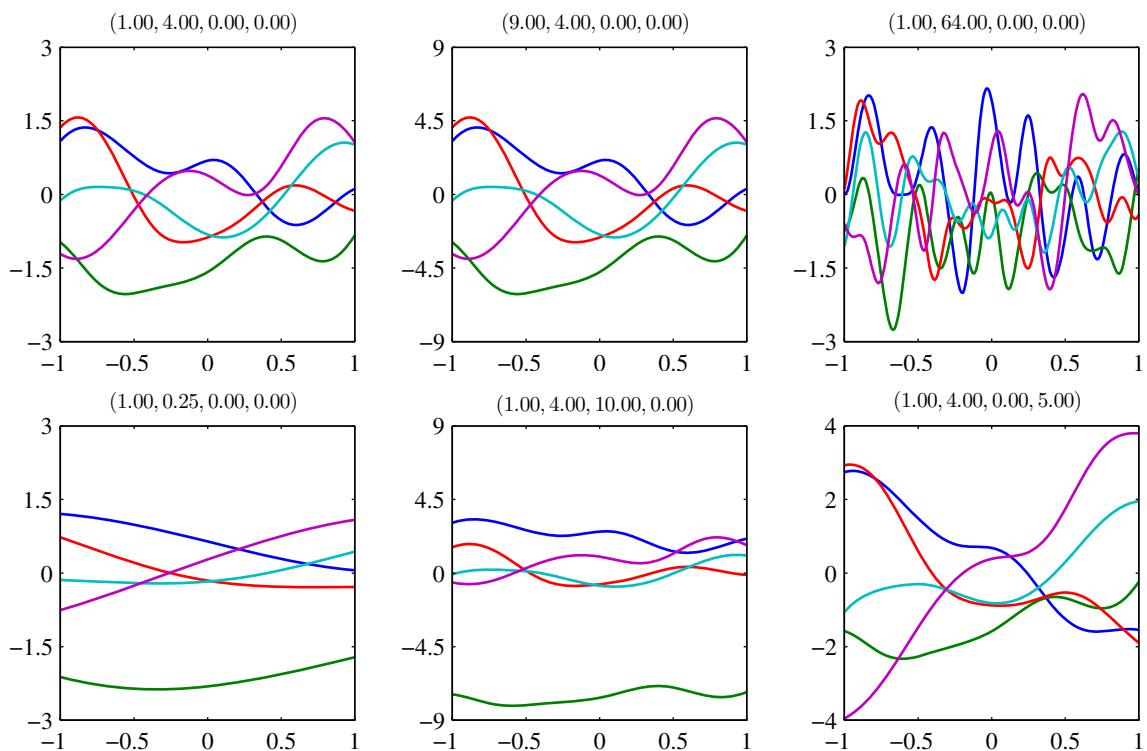


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$. Using the results (2.81) and (2.82), we see that the conditional distribution $p(t_{N+1} | \mathbf{t})$ is a Gaussian distribution with mean and covariance given by

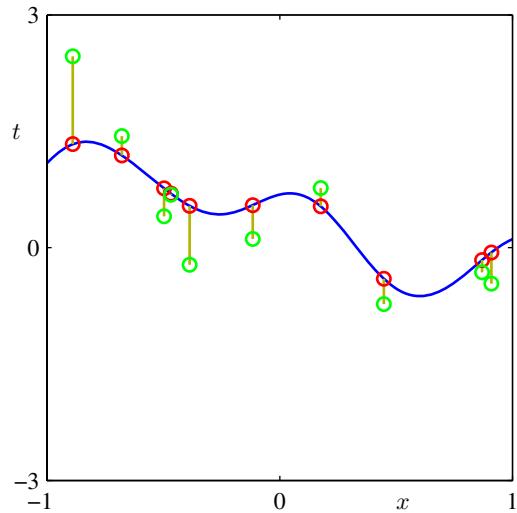
$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (6.66)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \quad (6.67)$$

These are the key results that define Gaussian process regression. Because the vector \mathbf{k} is a function of the test point input value \mathbf{x}_{N+1} , we see that the predictive distribution is a Gaussian whose mean and variance both depend on \mathbf{x}_{N+1} . An example of Gaussian process regression is shown in Figure 6.8.

任何合法的核函数都能使协差矩阵C也是合法的(也就是正定),因此对核函数的选择没有限制,只需要合法即可。

Figure 6.6 Illustration of the sampling of data points $\{t_n\}$ from a Gaussian process. The blue curve shows a sample function from the Gaussian process prior over functions, and the red points show the values of y_n obtained by evaluating the function at a set of input values $\{x_n\}$. The corresponding values of $\{t_n\}$, shown in green, are obtained by adding independent Gaussian noise to each of the $\{y_n\}$.



suitable kernels.

Note that the mean (6.66) of the predictive distribution can be written, as a function of \mathbf{x}_{N+1} , in the form

$$m(\mathbf{x}_{N+1}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}_{N+1}) \quad (6.68)$$

where a_n is the n^{th} component of $\mathbf{C}_N^{-1}\mathbf{t}$. Thus, if the kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$ depends only on the distance $\|\mathbf{x}_n - \mathbf{x}_m\|$, then we obtain an expansion in radial basis functions.

The results (6.66) and (6.67) define the predictive distribution for Gaussian process regression with an arbitrary kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$. In the particular case in which the kernel function $k(\mathbf{x}, \mathbf{x}')$ is defined in terms of a finite set of basis functions, we can derive the results obtained previously in Section 3.3.2 for linear regression starting from the Gaussian process viewpoint.

For such models, we can therefore obtain the predictive distribution either by taking a parameter space viewpoint and using the linear regression result or by taking a function space viewpoint and using the Gaussian process result.

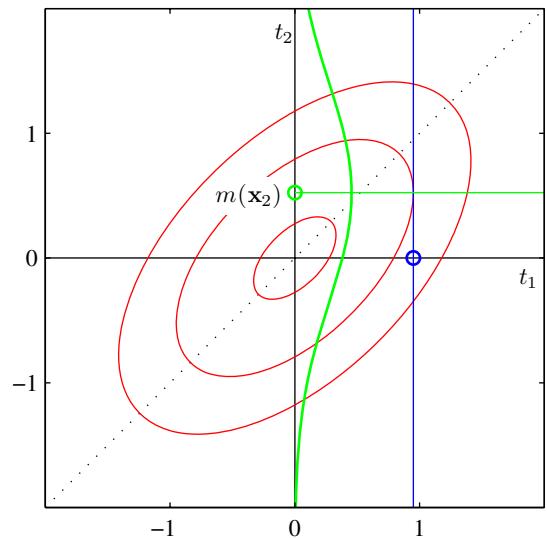
The central computational operation in using Gaussian processes will involve the inversion of a matrix of size $N \times N$, for which standard methods require $O(N^3)$ computations. By contrast, in the basis function model we have to invert a matrix \mathbf{S}_N of size $M \times M$, which has $O(M^3)$ computational complexity. Note that for both viewpoints, the matrix inversion must be performed once for the given training set. For each new test point, both methods require a vector-matrix multiply, which has cost $O(N^2)$ in the Gaussian process case and $O(M^2)$ for the linear basis function model. If the number M of basis functions is smaller than the number N of data points, it will be computationally more efficient to work in the basis function

GPR vs BLR优缺点比较

高斯过程回归与贝叶斯
线性回归的比较
当核函数对应的基函数
个数是有限时，可以高斯
过程回归推得贝叶斯线
性回归中间结果

Exercise 6.21

Figure 6.7 Illustration of the mechanism of Gaussian process regression for the case of one training point and one test point, in which the red ellipses show contours of the joint distribution $p(t_1, t_2)$. Here t_1 is the training data point, and conditioning on the value of t_1 , corresponding to the vertical blue line, we obtain $p(t_2|t_1)$ shown as a function of t_2 by the green curve.



framework. However, an advantage of a Gaussian processes viewpoint is that we can consider covariance functions that can only be expressed in terms of an infinite number of basis functions. []

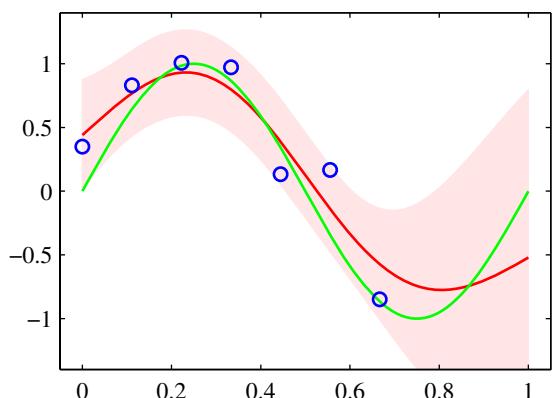
高斯过程的应用
与发展

For large training data sets, however, the direct application of Gaussian process methods can become infeasible, and so a range of approximation schemes have been developed that have better scaling with training set size than the exact approach (Gibbs, 1997; Tresp, 2001; Smola and Bartlett, 2001; Williams and Seeger, 2001; Csató and Opper, 2002; Seeger *et al.*, 2003). Practical issues in the application of Gaussian processes are discussed in Bishop and Nabney (2008).

We have introduced Gaussian process regression for the case of a single target variable. The extension of this formalism to multiple target variables, known as co-kriging (Cressie, 1993), is straightforward. Various other extensions of Gaus-

Exercise 6.23

Figure 6.8 Illustration of Gaussian process regression applied to the sinusoidal data set in Figure A.6 in which the three right-most data points have been omitted. The green curve shows the sinusoidal function from which the data points, shown in blue, are obtained by sampling and addition of Gaussian noise. The red line shows the mean of the Gaussian process predictive distribution, and the shaded region corresponds to plus and minus two standard deviations. Notice how the uncertainty increases in the region to the right of the data points.



sian process regression have also been considered, for purposes such as modelling the distribution over low-dimensional manifolds for unsupervised learning (Bishop *et al.*, 1998a) and the solution of stochastic differential equations (Graepel, 2003).]

6.4.3 Learning the hyperparameters

The predictions of a Gaussian process model will depend, in part, on the choice of covariance function. In practice, rather than fixing the covariance function, we may prefer to use a parametric family of functions and then infer the parameter values from the data. These parameters govern such things as the length scale of the correlations and the precision of the noise and correspond to the hyperparameters in a standard parametric model.

Techniques for learning the hyperparameters are based on the evaluation of the likelihood function $p(\mathbf{t}|\theta)$ where θ denotes the hyperparameters of the Gaussian process model. The simplest approach is to make a point estimate of θ by maximizing the log likelihood function. Because θ represents a set of hyperparameters for the regression problem, this can be viewed as analogous to the type 2 maximum likelihood procedure for linear regression models. Maximization of the log likelihood can be done using efficient gradient-based optimization algorithms such as conjugate gradients (Fletcher, 1987; Nocedal and Wright, 1999; Bishop and Nabney, 2008).

The log likelihood function for a Gaussian process regression model is easily evaluated using the standard form for a multivariate Gaussian distribution, giving

$$\ln p(\mathbf{t}|\theta) = -\frac{1}{2} \ln |\mathbf{C}_N| - \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2} \ln(2\pi). \quad (6.69)$$

$$= \int p(\mathbf{t}|\mathbf{w}) p(\mathbf{w}|\theta) d\mathbf{w}; \text{ 而 GPR 中 model evidence 可直接给出, 借用 P30 的第一句话}$$

For nonlinear optimization, we also need the gradient of the log likelihood function with respect to the parameter vector θ . We shall assume that evaluation of the derivatives of \mathbf{C}_N is straightforward, as would be the case for the covariance functions considered in this chapter. Making use of the result (C.21) for the derivative of \mathbf{C}_N^{-1} , together with the result (C.22) for the derivative of $\ln |\mathbf{C}_N|$, we obtain

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\theta) = -\frac{1}{2} \text{Tr} \left(\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \mathbf{C}_N^{-1} \mathbf{t}. \quad (6.70)$$

Because $\ln p(\mathbf{t}|\theta)$ will in general be a nonconvex function, it can have multiple maxima.

// It is straightforward to introduce a prior over θ and to maximize the log posterior using gradient-based methods. In a fully Bayesian treatment, we need to evaluate marginals over θ weighted by the product of the prior $p(\theta)$ and the likelihood function $p(\mathbf{t}|\theta)$. In general, however, exact marginalization will be intractable, and we must resort to approximations.]

[The Gaussian process regression model gives a predictive distribution whose mean and variance are functions of the input vector \mathbf{x} . However, we have assumed that the contribution to the predictive variance arising from the additive noise, governed by the parameter β , is a constant. For some problems, known as *heteroscedastic*, the noise variance itself will also depend on \mathbf{x} . To model this, we can extend the

GPR 的扩展

超参数 θ : 核函数中的超参数 (比如余弦或 RBF 形式)
而核函数时底方是 α) + β

也就是 max Model Evidence
事实上, 超参数还应包括基础

Section 3.5
数 (也就是核函数), 只不过
我们总是事先设定好, 而无
法像 α , β 一样进行自动
优化

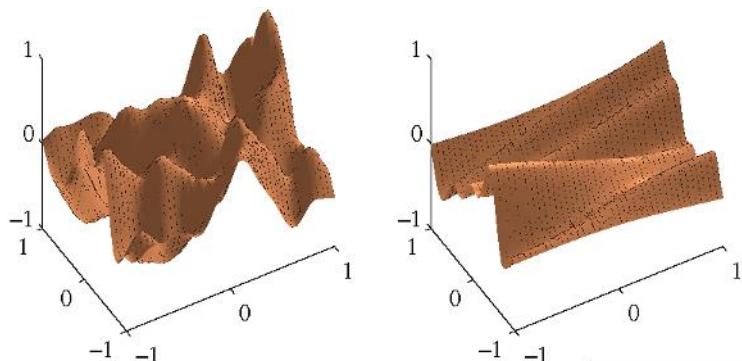
就是:
By working directly
with the covariance
function we have
implicitly marginalized
over the distribution of
weights.

即式 (6.53) 隐式地完成
了对 \mathbf{w} 的积分。

除了可以像上面那样选择
max evidence 来确定超参数
还可以引入贝叶斯先验, 进行
fully Bayesian treatment。

$P(\theta|\mathbf{t}) = \frac{P(\theta) P(\mathbf{t}|\theta)}{\int P(\theta) P(\mathbf{t}|\theta) d\theta}$,
若想计算 $P(\theta|\mathbf{t})$, 需计算
分子所求的积分, 一般
是 intractable, 需借助
近似算法。

Figure 6.9 Samples from the ARD prior for Gaussian processes, in which the kernel function is given by (6.71). The left plot corresponds to $\eta_1 = \eta_2 = 1$, and the right plot corresponds to $\eta_1 = 1, \eta_2 = 0.01$.



Gaussian process framework by introducing a second Gaussian process to represent the dependence of β on the input \mathbf{x} (Goldberg *et al.*, 1998). Because β is a variance, and hence nonnegative, we use the Gaussian process to model $\ln \beta(\mathbf{x})$. 33

6.4.4 Automatic relevance determination

在核函数中引入衰减参数，并和 β 一样通过最大似然确定，相当于基函数也是通过学习得到。(实际上， β 就是核函数中的一个参数)

多层神经网络最后的高层输出为视觉基函数，其与输出层的权重即为基函数的系数。可以看到，DNN 的基函数并不是固定的，也会有参数，可以通过学习得到。这一层的基函数的过程就是表示学习。

In the previous section, we saw how maximum likelihood could be used to determine a value for the correlation length-scale parameter in a Gaussian process. This technique can usefully be extended by incorporating a separate parameter for each input variable (Rasmussen and Williams, 2006). The result, as we shall see, is that the optimization of these parameters by maximum likelihood allows the relative importance of different inputs to be inferred from the data. This represents an example in the Gaussian process context of *automatic relevance determination*, or ARD, which was originally formulated in the framework of neural networks (MacKay, 1994; Neal, 1996). The mechanism by which appropriate inputs are preferred is discussed in Section 7.2.2. 7.2节介绍的模型也是上面想的一个例子，只不过7.2中并不是以高斯过程的形式，而是以参数模型的形式写出来的。7.2.1节给出了该模型通过起参数优化产生稀疏解的原因。

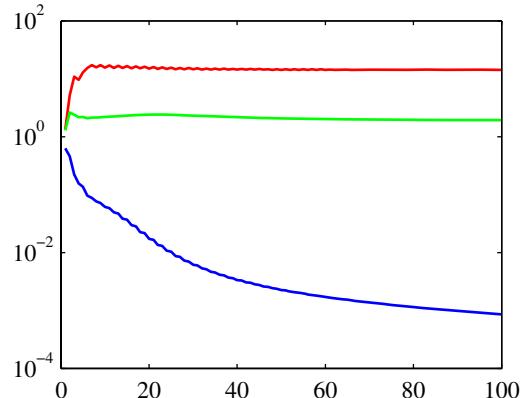
$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp \left\{ -\frac{1}{2} \sum_{i=1}^2 \eta_i (x_i - x'_i)^2 \right\}. \quad (6.71)$$

采用类似形式的核函数有利于获得稀疏解的好处

Samples from the resulting prior over functions $y(\mathbf{x})$ are shown for two different settings of the precision parameters η_i in Figure 6.9. We see that, as a particular parameter η_i becomes small, the function becomes relatively insensitive to the corresponding input variable x_i . By adapting these parameters to a data set using maximum likelihood, it becomes possible to detect input variables that have little effect on the predictive distribution, because the corresponding values of η_i will be small. This can be useful in practice because it allows such inputs to be discarded.

ARD is illustrated using a simple synthetic data set having three inputs x_1, x_2 and x_3 (Nabney, 2002) in Figure 6.10. The target variable t , is generated by sampling 100 values of x_1 from a Gaussian, evaluating the function $\sin(2\pi x_1)$, and then adding

Figure 6.10 Illustration of automatic relevance determination in a Gaussian process for a synthetic problem having three inputs x_1 , x_2 , and x_3 , for which the curves show the corresponding values of the hyperparameters η_1 (red), η_2 (green), and η_3 (blue) as a function of the number of iterations when optimizing the marginal likelihood. Details are given in the text. Note the logarithmic scale on the vertical axis.



Gaussian noise. Values of x_2 are given by copying the corresponding values of x_1 and adding noise, and values of x_3 are sampled from an independent Gaussian distribution. Thus x_1 is a good predictor of t , x_2 is a more noisy predictor of t , and x_3 has only chance correlations with t . The marginal likelihood for a Gaussian process with ARD parameters η_1, η_2, η_3 is optimized using the scaled conjugate gradients algorithm. We see from Figure 6.10 that η_1 converges to a relatively large value, η_2 converges to a much smaller value, and η_3 becomes very small indicating that x_3 is irrelevant for predicting t .

举一反三 [The ARD framework is easily incorporated into the exponential-quadratic kernel (6.63) to give the following form of kernel function, which has been found useful for applications of Gaussian processes to a range of regression problems]

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left\{ -\frac{1}{2} \sum_{i=1}^D \eta_i (x_{ni} - x_{mi})^2 \right\} + \theta_2 + \theta_3 \sum_{i=1}^D x_{ni} x_{mi} \quad (6.72)$$

where D is the dimensionality of the input space.

6.4.5 Gaussian processes for classification

In a probabilistic approach to classification, our goal is to model the posterior probabilities of the target variable for a new input vector, given a set of training data. These probabilities must lie in the interval $(0, 1)$, whereas a Gaussian process model makes predictions that lie on the entire real axis. [However, we can easily adapt Gaussian processes to classification problems by transforming the output of the Gaussian process using an appropriate nonlinear activation function.]

Consider first the two-class problem with a target variable $t \in \{0, 1\}$. If we define a Gaussian process over a function $a(\mathbf{x})$ and then transform the function using a logistic sigmoid $y = \sigma(a)$, given by (4.59), then we will obtain a non-Gaussian stochastic process over functions $y(\mathbf{x})$ where $y \in (0, 1)$. [This is illustrated for the case of a one-dimensional input space in Figure 6.11 in which the probability distri-

将 Gaussian processes 用于分类问题的处理方式

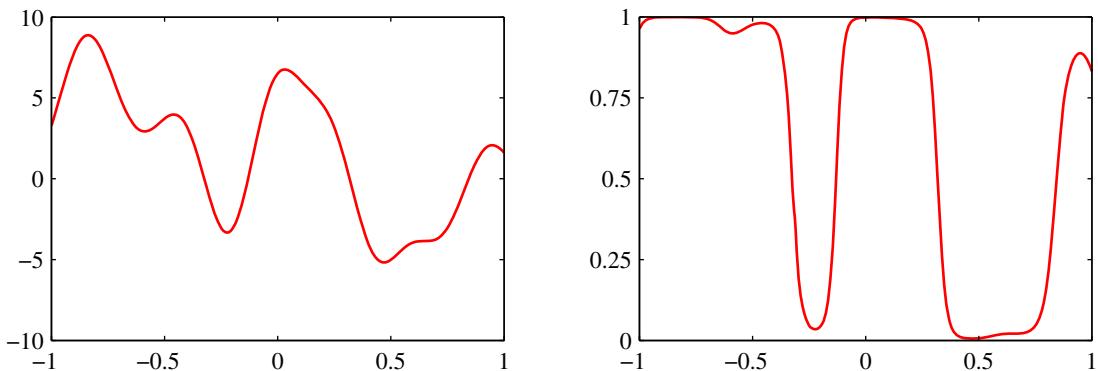


Figure 6.11 The left plot shows a sample from a Gaussian process prior over functions $a(\mathbf{x})$, and the right plot shows the result of transforming this sample using a logistic sigmoid function.

bution over the target variable t is then given by the Bernoulli distribution

$$p(t|a) = \sigma(a)^t(1 - \sigma(a))^{1-t}. \quad (6.73)$$

As usual, we denote the training set inputs by $\mathbf{x}_1, \dots, \mathbf{x}_N$ with corresponding observed target variables $\mathbf{t}_N = (t_1, \dots, t_N)^T$. We also consider a single test point \mathbf{x}_{N+1} with target value t_{N+1} . Our goal is to determine the predictive distribution $p(t_{N+1}|\mathbf{t})$, where we have left the conditioning on the input variables implicit. To do this we introduce a Gaussian process prior over the vector \mathbf{a}_{N+1} , which has components $a(\mathbf{x}_1), \dots, a(\mathbf{x}_{N+1})$. This in turn defines a non-Gaussian process over \mathbf{t}_{N+1} , and by conditioning on the training data \mathbf{t}_N we obtain the required predictive distribution. The Gaussian process prior for \mathbf{a}_{N+1} takes the form

$$p(\mathbf{a}_{N+1}) = \mathcal{N}(\mathbf{a}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1}). \quad (6.74)$$

校函数只是半正定
而协方差矩阵要
求正定，对应P308
最后一段的阐述

Unlike the regression case, the covariance matrix no longer includes a noise term because we assume that all of the training data points are correctly labelled. However, for numerical reasons it is convenient to introduce a noise-like term governed by a parameter ν that ensures that the covariance matrix is positive definite. Thus the covariance matrix \mathbf{C}_{N+1} has elements given by

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \nu \delta_{nm} \quad (6.75)$$

where $k(\mathbf{x}_n, \mathbf{x}_m)$ is any positive semidefinite kernel function of the kind considered in Section 6.2, and the value of ν is typically fixed in advance. We shall assume that the kernel function $k(\mathbf{x}, \mathbf{x}')$ is governed by a vector $\boldsymbol{\theta}$ of parameters, and we shall later discuss how $\boldsymbol{\theta}$ may be learned from the training data.

For two-class problems, it is sufficient to predict $p(t_{N+1} = 1 | \mathbf{t}_N)$ because the value of $p(t_{N+1} = 0 | \mathbf{t}_N)$ is then given by $1 - p(t_{N+1} = 1 | \mathbf{t}_N)$. The required

predictive distribution is given by

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1} \quad (6.76)$$

计算公式见式(6.71)

式(6.76)的近似计算：

法2：

对 $p(a_{N+1} | \mathbf{t}_N)$
进行高斯近似

Section 2.3
物理法支撑

法1：This integral is analytically intractable, and so may be approximated using sampling methods (Neal, 1997). Alternatively, we can consider techniques based on an analytical approximation. In Section 4.5.2, we derived the approximate formula (4.153) for the convolution of a logistic sigmoid with a Gaussian distribution. We can use this result to evaluate the integral in (6.76) provided we have a Gaussian approximation to the posterior distribution $p(a_{N+1} | \mathbf{t}_N)$. The usual justification for a Gaussian approximation to a posterior distribution is that the true posterior will tend to a Gaussian as the number of data points increases as a consequence of the central limit theorem. In the case of Gaussian processes, the number of variables grows with the number of data points, and so this argument does not apply directly. However, if we consider increasing the number of data points falling in a fixed region of \mathbf{x} space, then the corresponding uncertainty in the function $a(\mathbf{x})$ will decrease, again leading asymptotically to a Gaussian (Williams and Barber, 1998).

Three different approaches to obtaining a Gaussian approximation have been considered. One technique is based on *variational inference* (Gibbs and MacKay, 2000) and makes use of the local variational bound (10.144) on the logistic sigmoid. This allows the product of sigmoid functions to be approximated by a product of Gaussians thereby allowing the marginalization over \mathbf{a}_N to be performed analytically. The approach also yields a lower bound on the likelihood function $p(\mathbf{t}_N | \theta)$. The variational framework for Gaussian process classification can also be extended to multiclass ($K > 2$) problems by using a Gaussian approximation to the softmax function (Gibbs, 1997).

A second approach uses *expectation propagation* (Opper and Winther, 2000b; Minka, 2001b; Seeger, 2003). Because the true posterior distribution is unimodal, as we shall see shortly, the expectation propagation approach can give good results.

6.4.6 Laplace approximation for classification

The third approach to Gaussian process classification is based on the *Laplace approximation*, which we now consider in detail. In order to evaluate the predictive distribution (6.76), we seek a Gaussian approximation to the posterior distribution over a_{N+1} , which, using Bayes' theorem, is given by

$$\begin{aligned} p(a_{N+1} | \mathbf{t}_N) &= \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N \\ &= \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \end{aligned} \quad (6.77)$$

这个分布计算是个难题

利用Laplace近似进行计算

见下文

对 $p(a_{N+1} | \mathbf{t}_N)$ 进行高斯近似

Section 10.1

进行高斯近似

3种方法

第①种方法并不是基于式
(4.153)的结论,而是对 $p(a_{N+1})$
使用式(10.144)进行近似

Section 10.7

将式(6.76)的积分项完全
都变成高斯的,从而解
析计算

Section 4.4

$p(a_{N+1}|\mathbf{a}_N)$ 的计算：可解析结果

where we have used $p(\mathbf{t}_N|a_{N+1}, \mathbf{a}_N) = p(\mathbf{t}_N|\mathbf{a}_N)$. [The conditional distribution $p(a_{N+1}|\mathbf{a}_N)$ is obtained by invoking the results (6.66) and (6.67) for Gaussian process regression, to give

$$p(a_{N+1}|\mathbf{a}_N) = \mathcal{N}(a_{N+1}|\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}). \quad (6.78)$$

$p(a_N|\mathbf{t}_N)$ 的计算

基于 Laplace 近似

[We can therefore evaluate the integral in (6.77) by finding a Laplace approximation for the posterior distribution $p(\mathbf{a}_N|\mathbf{t}_N)$, and then using the standard result for the convolution of two Gaussian distributions.]

The prior $p(\mathbf{a}_N)$ is given by a zero-mean Gaussian process with covariance matrix \mathbf{C}_N , and the data term (assuming independence of the data points) is given by

$$p(\mathbf{t}_N|\mathbf{a}_N) = \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n). \quad (6.79)$$

We then obtain the Laplace approximation by Taylor expanding the logarithm of $p(\mathbf{a}_N|\mathbf{t}_N)$, which up to an additive normalization constant is given by the quantity

$$p(a_N|\mathbf{t}_N) = \frac{p(a_N)p(\mathbf{t}_N|a_N)}{p(\mathbf{t}_N)}$$

$$\propto p(a_N) p(\mathbf{t}_N|a_N)$$

$$\Psi(\mathbf{a}_N) = \ln p(\mathbf{a}_N) + \ln p(\mathbf{t}_N|\mathbf{a}_N) \\ = \left(-\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| \right) + (\mathbf{t}_N^T \mathbf{a}_N)$$

$$\text{分子 } p(\mathbf{t}_N|\theta) = \int p(a_N|\theta) p(\mathbf{t}_N|a_N) da_N$$

$$- \sum_{n=1}^N \ln(1 + e^{a_n}) + \text{const.} \quad (6.80)$$

Laplace approximation 是将 $p(a_N|\mathbf{t}_N)$ 近似为关于 a_N 的高斯分布，高斯分布的均值为 MAP，协差矩阵由取对数后 MAP 处的二阶导给定（实际上就是对 $p(a_N|\mathbf{t}_N)$ 在 MAP 处进行泰勒展开，并取二阶近似）。注意，用了近似的高斯分布的均值和方差并不是由 $p(a_N|\mathbf{t}_N)$ 的均值和方差给出。在计算过程中，分子 $p(\mathbf{t}_N|\theta)$ 为常数可忽略，也就是 $p(\mathbf{t}_N|\theta)$ 的值，即式(6.90)

First we need to find the mode of the posterior distribution, and this requires that we evaluate the gradient of $\Psi(\mathbf{a}_N)$, which is given by

$$\nabla \Psi(\mathbf{a}_N) = \mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N \quad (6.81)$$

where $\boldsymbol{\sigma}_N$ is a vector with elements $\sigma(a_n)$. We cannot simply find the mode by setting this gradient to zero, because $\boldsymbol{\sigma}_N$ depends nonlinearly on \mathbf{a}_N , and so we resort to an iterative scheme based on the Newton-Raphson method, which gives rise to an iterative reweighted least squares (IRLS) algorithm. This requires the second derivatives of $\Psi(\mathbf{a}_N)$, which we also require for the Laplace approximation anyway, and which are given by

$$\nabla \nabla \Psi(\mathbf{a}_N) = -\mathbf{W}_N - \mathbf{C}_N^{-1} \quad (6.82)$$

where \mathbf{W}_N is a diagonal matrix with elements $\sigma(a_n)(1 - \sigma(a_n))$, and we have used the result (4.88) for the derivative of the logistic sigmoid function. Note that these diagonal elements lie in the range $(0, 1/4)$, and hence \mathbf{W}_N is a positive definite matrix. Because \mathbf{C}_N (and hence its inverse) is positive definite by construction, and because the sum of two positive definite matrices is also positive definite, we see that the Hessian matrix $\mathbf{A} = -\nabla \nabla \Psi(\mathbf{a}_N)$ is positive definite and so the posterior distribution $p(\mathbf{a}_N|\mathbf{t}_N)$ is log convex and therefore has a single mode that is the global

Exercise 6.24
略，也就是说是说泰勒展开只针对 $\ln p(a_N|\mathbf{t}_N)$ 进行。完成对 $p(a_N|\mathbf{t}_N)$ 的计算和分得列而近似后，即为计算积分得列

maximum. The posterior distribution is not Gaussian, however, because the Hessian is a function of \mathbf{a}_N .

Exercise 6.25

Using the Newton-Raphson formula (4.92), the iterative update equation for \mathbf{a}_N is given by

$$\mathbf{a}_N^{\text{new}} = \mathbf{C}_N(\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} \{ \mathbf{t}_N - \boldsymbol{\sigma}_N + \mathbf{W}_N \mathbf{a}_N \}. \quad (6.83)$$

These equations are iterated until they converge to the mode which we denote by \mathbf{a}_N^* . At the mode, the gradient $\nabla \Psi(\mathbf{a}_N)$ will vanish, and hence \mathbf{a}_N^* will satisfy

$$\mathbf{a}_N^* = \mathbf{C}_N(\mathbf{t}_N - \boldsymbol{\sigma}_N). \quad \begin{matrix} \text{式(6.81)} \\ \text{注意, } \mathbf{W}_N \text{ evaluated at } \mathbf{a}_N^* \end{matrix} \quad (6.84)$$

Once we have found the mode \mathbf{a}_N^* of the posterior, we can evaluate the Hessian matrix given by

$$\mathbf{H} = -\nabla \nabla \Psi(\mathbf{a}_N) = \mathbf{W}_N + \mathbf{C}_N^{-1} \quad (6.85)$$

where the elements of \mathbf{W}_N are evaluated using \mathbf{a}_N^* . This defines our Gaussian approximation to the posterior distribution $p(\mathbf{a}_N | \mathbf{t}_N)$ given by

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N | \mathbf{a}_N^*, \mathbf{H}^{-1}). \quad (6.86)$$

练习(6.78)+(6.86)得

Exercise 6.26

[We can now combine this with (6.78) and hence evaluate the integral (6.77). Because this corresponds to a linear-Gaussian model, we can use the general result (2.115) to give

$$\mathbb{E}[a_{N+1} | \mathbf{t}_N] = \mathbf{k}^T(\mathbf{t}_N - \boldsymbol{\sigma}_N) \quad (6.87)$$

$$\text{var}[a_{N+1} | \mathbf{t}_N] = c - \mathbf{k}^T(\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1}\mathbf{k}. \quad (6.88)$$

Now that we have a Gaussian distribution for $p(a_{N+1} | \mathbf{t}_N)$, we can approximate the integral (6.76) using the result (4.153). As with the Bayesian logistic regression model of Section 4.5, if we are only interested in the decision boundary corresponding to $p(t_{N+1} | \mathbf{t}_N) = 0.5$, then we need only consider the mean and we can ignore the effect of the variance.]

[We also need to determine the parameters θ of the covariance function. One approach is to maximize the likelihood function given by $p(\mathbf{t}_N | \theta)$ for which we need expressions for the log likelihood and its gradient. If desired, suitable regularization terms can also be added, leading to a penalized maximum likelihood solution. The likelihood function is defined by

$$p(\mathbf{t}_N | \theta) = \int p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N | \theta) d\mathbf{a}_N. \quad \begin{matrix} \text{式(6.79)不是高斯分布} \\ \text{式(6.74)是高斯分布} \end{matrix} \quad (6.89)$$

This integral is analytically intractable, so again we make use of the Laplace approximation. Using the result (4.135), we obtain the following approximation for the log of the likelihood function

$$\ln p(\mathbf{t}_N | \theta) = \Psi(\mathbf{a}_N^*) - \frac{1}{2} \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2} \ln(2\pi) \quad (6.90)$$

勒近似, 参见前文 $p(\mathbf{a}_N | \mathbf{t}_N)$ 的 Laplace

近似部分。 $p(\mathbf{t}_N | \theta)$ 实际上就是

$$p(\mathbf{a}_N | \mathbf{t}_N) = \frac{p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N)}{p(\mathbf{t}_N)} \text{ in 分母}$$

算式(6.71), 即 $p(a_{N+1} | \mathbf{t}_N)$;
再代入式(6.76)中近似计算
 $p(t_{N+1} = 1 | \mathbf{t}_N)$, 至此, GPC
的预测部分结束

下面开始GPC起参数的
学习部分

注意, 不是对 $p(\mathbf{t}_N | \mathbf{a}_N)$ 进行
laplace近似, 得到一个近似
高斯分布再代入式(6.79)计
算两个高斯分布乘积的积
分。这里用 Laplace 近似针对
的是 $p(\mathbf{a}_N | \mathbf{t}_N)$, 等价地是取
 $\ln p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N | \theta)$ 的 = 分母

超参数 $\theta \rightarrow \mathbf{C}_N \rightarrow$ 高斯过程 $a_N \rightarrow t_N$

梯度计算
 θ 的变化会通过 C_N 影响
 $\text{MAP } a_N^*$

where $\Psi(\mathbf{a}_N^*) = \ln p(\mathbf{a}_N^* | \theta) + \ln p(\mathbf{t}_N | \mathbf{a}_N^*)$. We also need to evaluate the gradient of $\ln p(\mathbf{t}_N | \theta)$ with respect to the parameter vector θ . Note that changes in θ will cause changes in \mathbf{a}_N^* , leading to additional terms in the gradient. Thus, when we differentiate (6.90) with respect to θ , we obtain two sets of terms, the first arising from the dependence of the covariance matrix \mathbf{C}_N on θ , and the rest arising from dependence of \mathbf{a}_N^* on θ .

The terms arising from the explicit dependence on θ can be found by using (6.80) together with the results (C.21) and (C.22), and are given by

$$(C.21): \frac{\partial}{\partial \theta} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \theta} \mathbf{A}^{-1}$$

$$(C.22): \frac{\partial}{\partial \theta} \mathbf{I}_n |\mathbf{A}| = \text{Tr}(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \theta})$$

$$\mathbf{I}_n P(t_N | \theta) = -\frac{1}{2} \mathbf{a}_N^{*\top} \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{N}{2} \mathbf{I}_n(\mathbf{v}_2) - \frac{1}{2} \mathbf{I}_n |\mathbf{C}_N|$$

$$+ t_N^\top \mathbf{a}_N^* - \sum_{n=1}^N \mathbf{I}_n(1 + e^{a_n^*})$$

$$- \frac{1}{2} \mathbf{I}_n |\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{1}{2} \mathbf{I}_n(\mathbf{v}_2)$$

$$\begin{aligned} \frac{\partial \ln p(\mathbf{t}_N | \theta)}{\partial \theta_j} &= \frac{1}{2} \mathbf{a}_N^{*\top} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} \mathbf{C}_N^{-1} \mathbf{a}_N^* \\ &\quad - \frac{1}{2} \text{Tr} \left[(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{W}_N \frac{\partial \mathbf{C}_N}{\partial \theta_j} \right] \end{aligned} \quad (6.91)$$

To compute the terms arising from the dependence of \mathbf{a}_N^* on θ , we note that the Laplace approximation has been constructed such that $\Psi(\mathbf{a}_N)$ has zero gradient at $\mathbf{a}_N = \mathbf{a}_N^*$, and so $\Psi(\mathbf{a}_N^*)$ gives no contribution to the gradient as a result of its dependence on \mathbf{a}_N^* . This leaves the following contribution to the derivative with respect to a component θ_j of θ

将 $\partial \ln p(t_N | \theta)$ 的计算

分为两部分之和：

① 对中间变量 C_N 的偏导. $\frac{\partial C_N}{\partial \theta_j}$

② 对中间变量 a_N^* 的偏导. $\frac{\partial a_N^*}{\partial \theta_j}$

注意上式中 \mathbf{W}_N 为 a_N^* 的函数，在计算对 a_N^* 的偏导时不要漏掉了。

$$-\frac{1}{2} \sum_{n=1}^N \frac{\partial \ln |\mathbf{W}_N + \mathbf{C}_N^{-1}|}{\partial a_n^*} \frac{\partial a_n^*}{\partial \theta_j}$$

$$= -\frac{1}{2} \sum_{n=1}^N [(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{C}_N]_{nn} \sigma_n^*(1 - \sigma_n^*)(1 - 2\sigma_n^*) \frac{\partial a_n^*}{\partial \theta_j} \quad (6.92)$$

where $\sigma_n^* = \sigma(a_n^*)$, and again we have used the result (C.22) together with the definition of \mathbf{W}_N . We can evaluate the derivative of a_N^* with respect to θ_j by differentiating the relation (6.84) with respect to θ_j to give

$$\frac{\partial a_n^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N) - \mathbf{C}_N \mathbf{W}_N \frac{\partial a_n^*}{\partial \theta_j}. \quad (6.93)$$

Rearranging then gives

$$\frac{\partial a_n^*}{\partial \theta_j} = (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N). \quad (6.94)$$

Combining (6.91), (6.92), and (6.94), we can evaluate the gradient of the log likelihood function, which can be used with standard nonlinear optimization algorithms in order to determine a value for θ .

We can illustrate the application of the Laplace approximation for Gaussian processes using the synthetic two-class data set shown in Figure 6.12. Extension of the Laplace approximation to Gaussian processes involving $K > 2$ classes, using the softmax activation function, is straightforward (Williams and Barber, 1998).

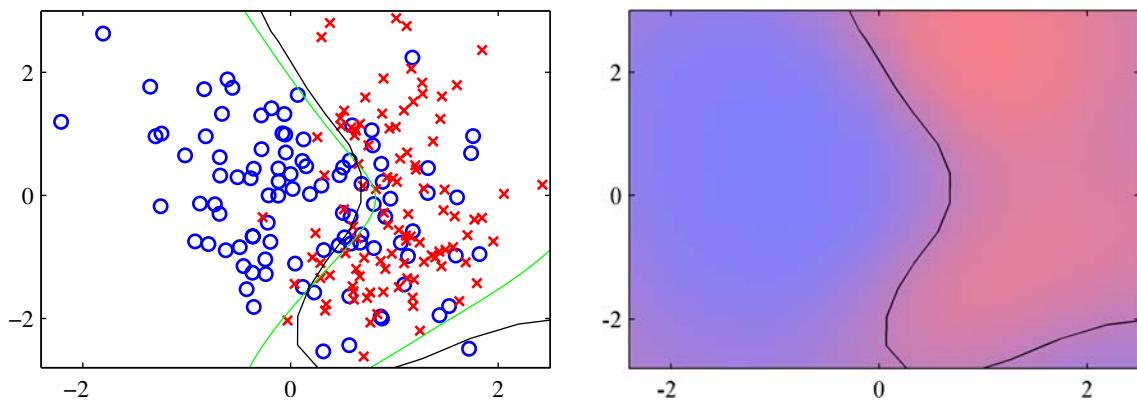


Figure 6.12 Illustration of the use of a Gaussian process for classification, showing the data on the left together with the optimal decision boundary from the true distribution in green, and the decision boundary from the Gaussian process classifier in black. On the right is the predicted posterior probability for the blue and red classes together with the Gaussian process decision boundary.

并不是在讨论GP与神经
网络的联系(相隔很远)
而是在讨论神
经网络函数的分布
是否是一个高斯过程,
其核函数有什么特点等

神经网络输出函数
的分布是高斯过程的情况

基于神经网络的高斯
过程的核函数的特点

6.4.7 Connection to neural networks

We have seen that the range of functions which can be represented by a neural network is governed by the number M of hidden units, and that, for sufficiently large M , a two-layer network can approximate any given function with arbitrary accuracy. In the framework of maximum likelihood, the number of hidden units needs to be limited (to a level dependent on the size of the training set) in order to avoid over-fitting. However, from a Bayesian perspective it makes little sense to limit the number of parameters in the network according to the size of the training set.

In a Bayesian neural network, the prior distribution over the parameter vector w , in conjunction with the network function $f(x, w)$, produces a prior distribution over functions from $y(x)$ where y is the vector of network outputs. Neal (1996) has shown that, for a broad class of prior distributions over w , the distribution of functions generated by a neural network will tend to a Gaussian process in the limit $M \rightarrow \infty$. It should be noted, however, that in this limit the output variables of the neural network become independent. One of the great merits of neural networks is that the outputs share the hidden units and so they can ‘borrow statistical strength’ from each other, that is, the weights associated with each hidden unit are influenced by all of the output variables not just by one of them. This property is therefore lost in the Gaussian process limit.

We have seen that a Gaussian process is determined by its covariance (kernel) function. Williams (1998) has given explicit forms for the covariance in the case of two specific choices for the hidden unit activation function (probit and Gaussian). These kernel functions $k(x, x')$ are nonstationary, i.e. they cannot be expressed as a function of the difference $x - x'$, as a consequence of the Gaussian weight prior being centred on zero which breaks translation invariance in weight space.

**核函数超参数
阅读**

By working directly with the covariance function we have implicitly marginalized over the distribution of weights. If the weight prior is governed by hyperparameters, then their values will determine the length scales of the distribution over functions, as can be understood by studying the examples in Figure 5.11 for the case of a finite number of hidden units. Note that we cannot marginalize out the hyperparameters analytically, and must instead resort to techniques of the kind discussed in Section 6.4.

Exercises

- 6.1** (**) **www** Consider the dual formulation of the least squares linear regression problem given in Section 6.1. Show that the solution for the components a_n of the vector \mathbf{a} can be expressed as a linear combination of the elements of the vector $\phi(\mathbf{x}_n)$. Denoting these coefficients by the vector \mathbf{w} , show that the dual of the dual formulation is given by the original representation in terms of the parameter vector \mathbf{w} .
- 6.2** (**) In this exercise, we develop a dual formulation of the perceptron learning algorithm. Using the perceptron learning rule (4.55), show that the learned weight vector \mathbf{w} can be written as a linear combination of the vectors $t_n \phi(\mathbf{x}_n)$ where $t_n \in \{-1, +1\}$. Denote the coefficients of this linear combination by α_n and derive a formulation of the perceptron learning algorithm, and the predictive function for the perceptron, in terms of the α_n . Show that the feature vector $\phi(\mathbf{x})$ enters only in the form of the kernel function $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$.
- 6.3** (*) The nearest-neighbour classifier (Section 2.5.2) assigns a new input vector \mathbf{x} to the same class as that of the nearest input vector \mathbf{x}_n from the training set, where in the simplest case, the distance is defined by the Euclidean metric $\|\mathbf{x} - \mathbf{x}_n\|^2$. By expressing this rule in terms of scalar products and then making use of kernel substitution, formulate the nearest-neighbour classifier for a general nonlinear kernel.
- 6.4** (*) In Appendix C, we give an example of a matrix that has positive elements but that has a negative eigenvalue and hence that is not positive definite. Find an example of the converse property, namely a 2×2 matrix with positive eigenvalues yet that has at least one negative element.
- 6.5** (*) **www** Verify the results (6.13) and (6.14) for constructing valid kernels.
- 6.6** (*) Verify the results (6.15) and (6.16) for constructing valid kernels.
- 6.7** (*) **www** Verify the results (6.17) and (6.18) for constructing valid kernels.
- 6.8** (*) Verify the results (6.19) and (6.20) for constructing valid kernels.
- 6.9** (*) Verify the results (6.21) and (6.22) for constructing valid kernels.
- 6.10** (*) Show that an excellent choice of kernel for learning a function $f(\mathbf{x})$ is given by $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}')$ by showing that a linear learning machine based on this kernel will always find a solution proportional to $f(\mathbf{x})$.

6.11 (★) By making use of the expansion (6.25), and then expanding the middle factor as a power series, show that the Gaussian kernel (6.23) can be expressed as the inner product of an infinite-dimensional feature vector.

6.12 (★★) **www** Consider the space of all possible subsets A of a given fixed set D . Show that the kernel function (6.27) corresponds to an inner product in a feature space of dimensionality $2^{|D|}$ defined by the mapping $\phi(A)$ where A is a subset of D and the element $\phi_U(A)$, indexed by the subset U , is given by

$$\phi_U(A) = \begin{cases} 1, & \text{if } U \subseteq A; \\ 0, & \text{otherwise.} \end{cases} \quad (6.95)$$

Here $U \subseteq A$ denotes that U is either a subset of A or is equal to A .

6.13 (★) Show that the Fisher kernel, defined by (6.33), remains invariant if we make a nonlinear transformation of the parameter vector $\theta \rightarrow \psi(\theta)$, where the function $\psi(\cdot)$ is invertible and differentiable.

6.14 (★) **www** Write down the form of the Fisher kernel, defined by (6.33), for the case of a distribution $p(\mathbf{x}|\boldsymbol{\mu}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{S})$ that is Gaussian with mean $\boldsymbol{\mu}$ and fixed covariance \mathbf{S} .

6.15 (★) By considering the determinant of a 2×2 Gram matrix, show that a positive-definite kernel function $k(x, x')$ satisfies the Cauchy-Schwartz inequality

$$k(x_1, x_2)^2 \leq k(x_1, x_1)k(x_2, x_2). \quad (6.96)$$

6.16 (★★) Consider a parametric model governed by the parameter vector \mathbf{w} together with a data set of input values $\mathbf{x}_1, \dots, \mathbf{x}_N$ and a nonlinear feature mapping $\phi(\mathbf{x})$. Suppose that the dependence of the error function on \mathbf{w} takes the form

$$J(\mathbf{w}) = f(\mathbf{w}^T \phi(\mathbf{x}_1), \dots, \mathbf{w}^T \phi(\mathbf{x}_N)) + g(\mathbf{w}^T \mathbf{w}) \quad (6.97)$$

where $g(\cdot)$ is a monotonically increasing function. By writing \mathbf{w} in the form

$$\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n) + \mathbf{w}_\perp \quad (6.98)$$

where $\mathbf{w}_\perp^T \phi(\mathbf{x}_n) = 0$ for all n ,

✓ show that the value of \mathbf{w} that minimizes $J(\mathbf{w})$ takes the form of a linear combination of the basis functions $\phi(\mathbf{x}_n)$ for $n = 1, \dots, N$.

6.17 (★★) **www** Consider the sum-of-squares error function (6.39) for data having noisy inputs, where $\nu(\xi)$ is the distribution of the noise. Use the calculus of variations to minimize this error function with respect to the function $y(\mathbf{x})$, and hence show that the optimal solution is given by an expansion of the form (6.40) in which the basis functions are given by (6.41).

6.18 (*) Consider a Nadaraya-Watson model with one input variable x and one target variable t having Gaussian components with isotropic covariances, so that the covariance matrix is given by $\sigma^2 \mathbf{I}$ where \mathbf{I} is the unit matrix. Write down expressions for the conditional density $p(t|x)$ and for the conditional mean $\mathbb{E}[t|x]$ and variance $\text{var}[t|x]$, in terms of the kernel function $k(x, x_n)$.

6.19 (**) Another viewpoint on kernel regression comes from a consideration of regression problems in which the input variables as well as the target variables are corrupted with additive noise. Suppose each target value t_n is generated as usual by taking a function $y(\mathbf{z}_n)$ evaluated at a point \mathbf{z}_n , and adding Gaussian noise. The value of \mathbf{z}_n is not directly observed, however, but only a noise corrupted version $\mathbf{x}_n = \mathbf{z}_n + \boldsymbol{\xi}_n$ where the random variable $\boldsymbol{\xi}$ is governed by some distribution $g(\boldsymbol{\xi})$. Consider a set of observations $\{\mathbf{x}_n, t_n\}$, where $n = 1, \dots, N$, together with a corresponding sum-of-squares error function defined by averaging over the distribution of input noise to give

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n - \boldsymbol{\xi}_n) - t_n\}^2 g(\boldsymbol{\xi}_n) d\boldsymbol{\xi}_n. \quad (6.99)$$

By minimizing E with respect to the function $y(\mathbf{z})$ using the calculus of variations (Appendix D), show that optimal solution for $y(\mathbf{x})$ is given by a Nadaraya-Watson kernel regression solution of the form (6.45) with a kernel of the form (6.46).

6.20 (**) **www** Verify the results (6.66) and (6.67).

6.21 (**) **www** Consider a Gaussian process regression model in which the kernel function is defined in terms of a fixed set of nonlinear basis functions. Show that the predictive distribution is identical to the result (3.58) obtained in Section 3.3.2 for the Bayesian linear regression model. To do this, note that both models have Gaussian predictive distributions, and so it is only necessary to show that the conditional mean and variance are the same. For the mean, make use of the matrix identity (C.6), and for the variance, make use of the matrix identity (C.7).

6.22 (**) Consider a regression problem with N training set input vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ and L test set input vectors $\mathbf{x}_{N+1}, \dots, \mathbf{x}_{N+L}$, and suppose we define a Gaussian process prior over functions $t(\mathbf{x})$. Derive an expression for the joint predictive distribution for $t(\mathbf{x}_{N+1}), \dots, t(\mathbf{x}_{N+L})$, given the values of $t(\mathbf{x}_1), \dots, t(\mathbf{x}_N)$. Show the marginal of this distribution for one of the test observations t_j where $N+1 \leq j \leq N+L$ is given by the usual Gaussian process regression result (6.66) and (6.67).

6.23 (**) **www** Consider a Gaussian process regression model in which the target variable \mathbf{t} has dimensionality D . Write down the conditional distribution of \mathbf{t}_{N+1} for a test input vector \mathbf{x}_{N+1} , given a training set of input vectors $\mathbf{x}_1, \dots, \mathbf{x}_{N+L}$ and corresponding target observations $\mathbf{t}_1, \dots, \mathbf{t}_N$.

6.24 (*) Show that a diagonal matrix \mathbf{W} whose elements satisfy $0 < W_{ii} < 1$ is positive definite. Show that the sum of two positive definite matrices is itself positive definite.

- 6.25** (★) **www** Using the Newton-Raphson formula (4.92), derive the iterative update formula (6.83) for finding the mode \mathbf{a}_N^* of the posterior distribution in the Gaussian process classification model.
- 6.26** (★) Using the result (2.115), derive the expressions (6.87) and (6.88) for the mean and variance of the posterior distribution $p(a_{N+1}|\mathbf{t}_N)$ in the Gaussian process classification model.
- 6.27** (★★★) Derive the result (6.90) for the log likelihood function in the Laplace approximation framework for Gaussian process classification. Similarly, derive the results (6.91), (6.92), and (6.94) for the terms in the gradient of the log likelihood.

7

Sparse Kernel Machines

总结上一章并提出
本章的课题

[In the previous chapter, we explored a variety of learning algorithms based on non-linear kernels. One of the significant limitations of many such algorithms is that the kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$ must be evaluated for all possible pairs \mathbf{x}_n and \mathbf{x}_m of training points, which can be computationally infeasible during training and can lead to excessive computation times when making predictions for new data points.]

// In this chapter we shall look at kernel-based algorithms that have *sparse* solutions, so that predictions for new inputs depend only on the kernel function evaluated at a subset of the training data points.]

We begin by looking in some detail at the *support vector machine* (SVM), which became popular in some years ago for solving problems in classification, regression, and novelty detection. An important property of support vector machines is that the determination of the model parameters corresponds to a convex optimization problem, and so any local solution is also a global optimum. Because the discussion of support vector machines makes extensive use of Lagrange multipliers, the reader is

encouraged to review the key concepts covered in Appendix E. Additional information on support vector machines can be found in Vapnik (1995), Burges (1998), Cristianini and Shawe-Taylor (2000), Müller *et al.* (2001), Schölkopf and Smola (2002), and Herbrich (2002).

Section 7.2

The SVM is a decision machine and so does not provide posterior probabilities. We have already discussed some of the benefits of determining probabilities in Section 1.5.4. An alternative sparse kernel technique, known as the *relevance vector machine* (RVM), is based on a Bayesian formulation and provides posterior probabilistic outputs, as well as having typically much sparser solutions than the SVM.

7.1. Maximum Margin Classifiers

We begin our discussion of support vector machines by returning to the two-class classification problem using linear models of the form

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b \quad (7.1)$$

where $\phi(\mathbf{x})$ denotes a fixed feature-space transformation, and we have made the bias parameter b explicit. Note that we shall shortly introduce a dual representation expressed in terms of kernel functions, which avoids having to work explicitly in feature space. The training data set comprises N input vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$, with corresponding target values t_1, \dots, t_N where $t_n \in \{-1, 1\}$, and new data points \mathbf{x} are classified according to the sign of $y(\mathbf{x})$.

We shall assume for the moment that the training data set is linearly separable in feature space, so that by definition there exists at least one choice of the parameters \mathbf{w} and b such that a function of the form (7.1) satisfies $y(\mathbf{x}_n) > 0$ for points having $t_n = +1$ and $y(\mathbf{x}_n) < 0$ for points having $t_n = -1$, so that $t_n y(\mathbf{x}_n) > 0$ for all training data points.

There may of course exist many such solutions that separate the classes exactly. In Section 4.1.7, we described the perceptron algorithm that is guaranteed to find a solution in a finite number of steps. The solution that it finds, however, will be dependent on the (arbitrary) initial values chosen for \mathbf{w} and b as well as on the order in which the data points are presented. If there are multiple solutions all of which classify the training data set exactly, then we should try to find the one that will give the smallest generalization error. The support vector machine approaches this problem through the concept of the *margin*, which is defined to be the smallest distance between the decision boundary and any of the samples, as illustrated in Figure 7.1.

Section 7.1.5

In support vector machines the decision boundary is chosen to be the one for which the margin is maximized. The maximum margin solution can be motivated using computational learning theory, also known as statistical learning theory. However, a simple insight into the origins of maximum margin has been given by Tong and Koller (2000) who consider a framework for classification based on a hybrid of generative and discriminative approaches. They first model the distribution over input vectors \mathbf{x} for each class using a Parzen density estimator with Gaussian kernels

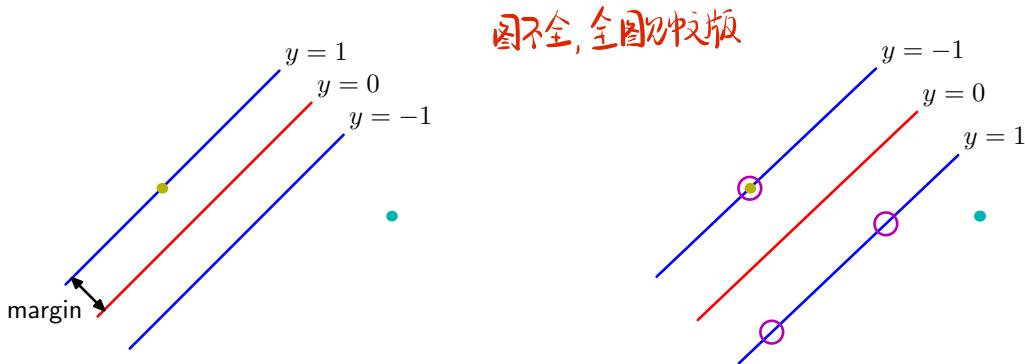


Figure 7.1 The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined by a subset of the data points, known as support vectors, which are indicated by the circles.

having a common parameter σ^2 . Together with the class priors, this defines an optimal misclassification-rate decision boundary. However, instead of using this optimal boundary, they determine the best hyperplane by minimizing the probability of error relative to the learned density model. In the limit $\sigma^2 \rightarrow 0$, the optimal hyperplane is shown to be the one having maximum margin. The intuition behind this result is that as σ^2 is reduced, the hyperplane is increasingly dominated by nearby data points relative to more distant ones. In the limit, the hyperplane becomes independent of data points that are not support vectors.

We shall see in Figure 10.13 that marginalization with respect to the prior distribution of the parameters in a Bayesian approach for a simple linearly separable data set leads to a decision boundary that lies in the middle of the region separating the data points. The large margin solution has similar behaviour.

Recall from Figure 4.1 that the perpendicular distance of a point \mathbf{x} from a hyperplane defined by $y(\mathbf{x}) = 0$ where $y(\mathbf{x})$ takes the form (7.1) is given by $|y(\mathbf{x})|/\|\mathbf{w}\|$. Furthermore, we are only interested in solutions for which all data points are correctly classified, so that $t_n y(\mathbf{x}_n) > 0$ for all n . Thus the distance of a point \mathbf{x}_n to the decision surface is given by

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}. \quad (7.2)$$

The margin is given by the perpendicular distance to the closest point \mathbf{x}_n from the data set, and we wish to optimize the parameters \mathbf{w} and b in order to maximize this distance. Thus the maximum margin solution is found by solving

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\} \quad (7.3)$$

where we have taken the factor $1/\|\mathbf{w}\|$ outside the optimization over n because \mathbf{w}

does not depend on n . Direct solution of this optimization problem would be very complex, and so we shall convert it into an equivalent problem that is much easier to solve. To do this we note that if we make the rescaling $\mathbf{w} \rightarrow \kappa\mathbf{w}$ and $b \rightarrow \kappa b$, then the distance from any point \mathbf{x}_n to the decision surface, given by $t_n y(\mathbf{x}_n)/\|\mathbf{w}\|$, is unchanged. We can use this freedom to set

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1 \quad (7.4)$$

for the point that is closest to the surface. In this case, all data points will satisfy the constraints

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N. \quad (7.5)$$

This is known as the canonical representation of the decision hyperplane. In the case of data points for which the equality holds, the constraints are said to be *active*, whereas for the remainder they are said to be *inactive*. By definition, there will always be at least one active constraint, because there will always be a closest point, and once the margin has been maximized there will be at least two active constraints. The optimization problem then simply requires that we maximize $\|\mathbf{w}\|^{-1}$, which is equivalent to minimizing $\|\mathbf{w}\|^2$, and so we have to solve the optimization problem

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.6)$$

subject to the constraints given by (7.5). The factor of $1/2$ in (7.6) is included for later convenience. This is an example of a *quadratic programming* problem in which we are trying to minimize a quadratic function subject to a set of linear inequality constraints. It appears that the bias parameter b has disappeared from the optimization. However, it is determined implicitly via the constraints, because these require that changes to $\|\mathbf{w}\|$ be compensated by changes to b . We shall see how this works shortly.

In order to solve this constrained optimization problem, we introduce Lagrange multipliers $a_n \geq 0$, with one multiplier a_n for each of the constraints in (7.5), giving the Lagrangian function

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{ t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 \} \quad (7.7)$$

where $\mathbf{a} = (a_1, \dots, a_N)^T$. Note the minus sign in front of the Lagrange multiplier term, because we are minimizing with respect to \mathbf{w} and b , and maximizing with respect to \mathbf{a} . Setting the derivatives of $L(\mathbf{w}, b, \mathbf{a})$ with respect to \mathbf{w} and b equal to zero, we obtain the following two conditions

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (7.8)$$

$$0 = \sum_{n=1}^N a_n t_n. \quad (7.9)$$

Eliminating \mathbf{w} and b from $L(\mathbf{w}, b, \mathbf{a})$ using these conditions then gives the *dual representation* of the maximum margin problem in which we maximize

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (7.10)$$

with respect to \mathbf{a} subject to the constraints

$$a_n \geq 0, \quad n = 1, \dots, N, \quad (7.11)$$

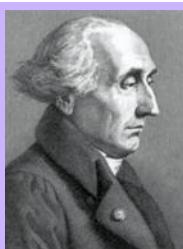
$$\sum_{n=1}^N a_n t_n = 0. \quad (7.12)$$

Here the kernel function is defined by $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$. Again, this takes the form of a quadratic programming problem in which we optimize a quadratic function of \mathbf{a} subject to a set of inequality constraints. We shall discuss techniques for solving such quadratic programming problems in Section 7.1.1.

The solution to a quadratic programming problem in M variables in general has computational complexity that is $O(M^3)$. In going to the dual formulation we have turned the original optimization problem, which involved minimizing (7.6) over M variables, into the dual problem (7.10), which has N variables. For a fixed set of basis functions whose number M is smaller than the number N of data points, the move to the dual problem appears disadvantageous. However, it allows the model to be reformulated using kernels, and so the maximum margin classifier can be applied efficiently to feature spaces whose dimensionality exceeds the number of data points, including infinite feature spaces. The kernel formulation also makes clear the role of the constraint that the kernel function $k(\mathbf{x}, \mathbf{x}')$ be positive definite, because this ensures that the Lagrangian function $\tilde{L}(\mathbf{a})$ is bounded below, giving rise to a well-defined optimization problem.

In order to classify new data points using the trained model, we evaluate the sign of $y(\mathbf{x})$ defined by (7.1). This can be expressed in terms of the parameters $\{a_n\}$ and the kernel function by substituting for \mathbf{w} using (7.8) to give

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b. \quad (7.13)$$



Joseph-Louis Lagrange
1736–1813

Although widely considered to be a French mathematician, Lagrange was born in Turin in Italy. By the age of nineteen, he had already made important contributions to mathematics and had been appointed as Professor at the Royal Artillery School in Turin. For many

years, Euler worked hard to persuade Lagrange to move to Berlin, which he eventually did in 1766 where he succeeded Euler as Director of Mathematics at the Berlin Academy. Later he moved to Paris, narrowly escaping with his life during the French revolution thanks to the personal intervention of Lavoisier (the French chemist who discovered oxygen) who himself was later executed at the guillotine. Lagrange made key contributions to the calculus of variations and the foundations of dynamics.

In Appendix E, we show that a constrained optimization of this form satisfies the **Karush-Kuhn-Tucker (KKT) conditions**, which in this case require that the following three properties hold

$$a_n \geq 0 \quad (7.14)$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0 \quad (7.15)$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0. \quad (7.16)$$

Thus for every data point, either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1$. Any data point for which $a_n = 0$ will not appear in the sum in (7.13) and hence plays no role in making predictions for new data points. The remaining data points are called *support vectors*, and because they satisfy $t_n y(\mathbf{x}_n) = 1$, they correspond to points that lie on the maximum margin hyperplanes in feature space, as illustrated in Figure 7.1. This property is central to the practical applicability of support vector machines. Once the model is trained, a significant proportion of the data points can be discarded and only the support vectors retained.

Having solved the quadratic programming problem and found a value for \mathbf{a} , we can then determine the value of the threshold parameter b by noting that any support vector \mathbf{x}_n satisfies $t_n y(\mathbf{x}_n) = 1$. Using (7.13) this gives

$$t_n \left(\sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1 \quad (7.17)$$

where \mathcal{S} denotes the set of indices of the support vectors. Although we can solve this equation for b using an arbitrarily chosen support vector \mathbf{x}_n , a numerically more stable solution is obtained by first multiplying through by t_n , making use of $t_n^2 = 1$, and then averaging these equations over all support vectors and solving for b to give

$$b = \frac{1}{N_{\mathcal{S}}} \sum_{n \in \mathcal{S}} \left(t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right) \quad (7.18)$$

where $N_{\mathcal{S}}$ is the total number of support vectors.

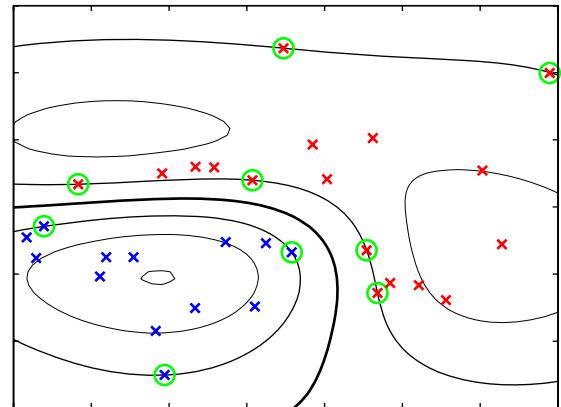
For later comparison with alternative models, we can express the maximum-margin classifier in terms of the minimization of an error function, with a simple quadratic regularizer, in the form

$$\sum_{n=1}^N E_\infty(y(\mathbf{x}_n)t_n - 1) + \lambda \|\mathbf{w}\|^2 \quad (7.19)$$

where $E_\infty(z)$ is a function that is zero if $z \geq 0$ and ∞ otherwise and ensures that the constraints (7.5) are satisfied. Note that as long as the regularization parameter satisfies $\lambda > 0$, its precise value plays no role.

Figure 7.2 shows an example of the classification resulting from training a support vector machine on a simple synthetic data set using a Gaussian kernel of the

Figure 7.2 Example of synthetic data from two classes in two dimensions showing contours of constant $y(\mathbf{x})$ obtained from a support vector machine having a Gaussian kernel function. Also shown are the decision boundary, the margin boundaries, and the support vectors.



form (6.23). Although the data set is not linearly separable in the two-dimensional data space \mathbf{x} , it is linearly separable in the nonlinear feature space defined implicitly by the nonlinear kernel function. Thus the training data points are perfectly separated in the original data space.

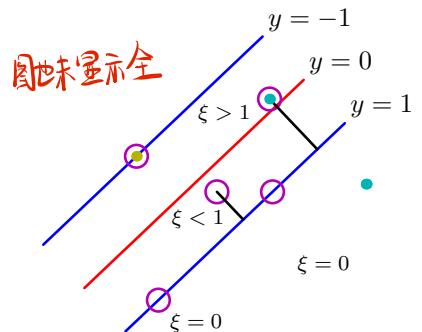
This example also provides a geometrical insight into the origin of sparsity in the SVM. The maximum margin hyperplane is defined by the location of the support vectors. Other data points can be moved around freely (so long as they remain outside the margin region) without changing the decision boundary, and so the solution will be independent of such data points.

7.1.1 Overlapping class distributions

So far, we have assumed that the training data points are linearly separable in the feature space $\phi(\mathbf{x})$. The resulting support vector machine will give exact separation of the training data in the original input space \mathbf{x} , although the corresponding decision boundary will be nonlinear. In practice, however, the class-conditional distributions may overlap, in which case exact separation of the training data can lead to poor generalization.

We therefore need a way to modify the support vector machine so as to allow some of the training points to be misclassified. From (7.19) we see that in the case of separable classes, we implicitly used an error function that gave infinite error if a data point was misclassified and zero error if it was classified correctly, and then optimized the model parameters to maximize the margin. We now modify this approach so that data points are allowed to be on the ‘wrong side’ of the margin boundary, but with a penalty that increases with the distance from that boundary. For the subsequent optimization problem, it is convenient to make this penalty a linear function of this distance. To do this, we introduce *slack variables*, $\xi_n \geq 0$ where $n = 1, \dots, N$, with one slack variable for each training data point (Bennett, 1992; Cortes and Vapnik, 1995). These are defined by $\xi_n = 0$ for data points that are on or inside the correct margin boundary and $\xi_n = |t_n - y(\mathbf{x}_n)|$ for other points. Thus a data point that is on the decision boundary $y(\mathbf{x}_n) = 0$ will have $\xi_n = 1$, and points

Figure 7.3 Illustration of the slack variables $\xi_n \geq 0$. Data points with circles around them are support vectors.



with $\xi_n > 1$ will be misclassified. The exact classification constraints (7.5) are then replaced with

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, \quad n = 1, \dots, N \quad (7.20)$$

in which the slack variables are constrained to satisfy $\xi_n \geq 0$. Data points for which $\xi_n = 0$ are correctly classified and are either on the margin or on the correct side of the margin. Points for which $0 < \xi_n \leq 1$ lie inside the margin, but on the correct side of the decision boundary, and those data points for which $\xi_n > 1$ lie on the wrong side of the decision boundary and are misclassified, as illustrated in Figure 7.3. This is sometimes described as relaxing the hard margin constraint to give a *soft margin* and allows some of the training set data points to be misclassified. Note that while slack variables allow for overlapping class distributions, this framework is still sensitive to outliers because the penalty for misclassification increases linearly with ξ .

Our goal is now to maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary. We therefore minimize

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.21)$$

where the parameter $C > 0$ controls the trade-off between the slack variable penalty and the margin. Because any point that is misclassified has $\xi_n > 1$, it follows that $\sum_n \xi_n$ is an upper bound on the number of misclassified points. The parameter C is therefore analogous to (the inverse of) a regularization coefficient because it controls the trade-off between minimizing training errors and controlling model complexity. In the limit $C \rightarrow \infty$, we will recover the earlier support vector machine for separable data.

We now wish to minimize (7.21) subject to the constraints (7.20) together with $\xi_n \geq 0$. The corresponding Lagrangian is given by

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n \quad (7.22)$$

where $\{a_n \geq 0\}$ and $\{\mu_n \geq 0\}$ are Lagrange multipliers. The corresponding set of KKT conditions are given by

$$a_n \geq 0 \quad (7.23)$$

$$t_n y(\mathbf{x}_n) - 1 + \xi_n \geq 0 \quad (7.24)$$

$$a_n (t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0 \quad (7.25)$$

$$\mu_n \geq 0 \quad (7.26)$$

$$\xi_n \geq 0 \quad (7.27)$$

$$\mu_n \xi_n = 0 \quad (7.28)$$

where $n = 1, \dots, N$.

We now optimize out \mathbf{w} , b , and $\{\xi_n\}$ making use of the definition (7.1) of $y(\mathbf{x})$ to give

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (7.29)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0 \quad (7.30)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n. \quad (7.31)$$

Using these results to eliminate \mathbf{w} , b , and $\{\xi_n\}$ from the Lagrangian, we obtain the dual Lagrangian in the form

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (7.32)$$

which is identical to the separable case, except that the constraints are somewhat different. To see what these constraints are, we note that $a_n \geq 0$ is required because these are Lagrange multipliers. Furthermore, (7.31) together with $\mu_n \geq 0$ implies $a_n \leq C$. We therefore have to minimize (7.32) with respect to the dual variables $\{a_n\}$ subject to ~~maximize~~

$$0 \leq a_n \leq C \quad (7.33)$$

$$\sum_{n=1}^N a_n t_n = 0 \quad (7.34)$$

for $n = 1, \dots, N$, where (7.33) are known as *box constraints*. This again represents a quadratic programming problem. If we substitute (7.29) into (7.1), we see that predictions for new data points are again made by using (7.13).

对支持向量的解读 We can now interpret the resulting solution. As before, a subset of the data points may have $a_n = 0$, in which case they do not contribute to the predictive

model (7.13). The remaining data points constitute the support vectors. These have $a_n > 0$ and hence from (7.25) must satisfy

$$t_n y(\mathbf{x}_n) = 1 - \xi_n. \quad (7.35)$$

If $a_n < C$, then (7.31) implies that $\mu_n > 0$, which from (7.28) requires $\xi_n = 0$ and hence such points lie on the margin. Points with $a_n = C$ can lie inside the margin and can either be correctly classified if $\xi_n \leq 1$ or misclassified if $\xi_n > 1$.

偏置b的计算 To determine the parameter b in (7.1), we note that those support vectors for which $0 < a_n < C$ have $\xi_n = 0$ so that $t_n y(\mathbf{x}_n) = 1$ and hence will satisfy

$$t_n \left(\sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1. \quad (7.36)$$

Again, a numerically stable solution is obtained by averaging to give

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left(t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right) \quad (7.37)$$

where \mathcal{M} denotes the set of indices of data points having $0 < a_n < C$.

SVM的另一种形式 ν -SVM An alternative, equivalent formulation of the support vector machine, known as the ν -SVM, has been proposed by Schölkopf *et al.* (2000). This involves maximizing the margin while minimizing the error. The objective function is

$$\tilde{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (7.38)$$

subject to the constraints

$$0 \leq a_n \leq 1/N \quad (7.39)$$

$$\sum_{n=1}^N a_n t_n = 0 \quad (7.40)$$

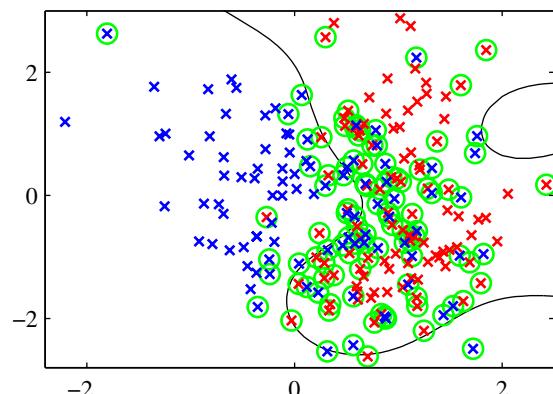
$$\sum_{n=1}^N a_n \geq \nu. \quad (7.41)$$

This approach has the advantage that the parameter ν , which replaces C , can be interpreted as both an upper bound on the fraction of *margin errors* (points for which $\xi_n > 0$ and hence which lie on the wrong side of the margin boundary and which may or may not be misclassified) and a lower bound on the fraction of support vectors. An example of the ν -SVM applied to a synthetic data set is shown in Figure 7.4. Here Gaussian kernels of the form $\exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$ have been used, with $\gamma = 0.45$.

Although predictions for new inputs are made using only the support vectors, the training phase (i.e., the determination of the parameters \mathbf{a} and b) makes use of the whole data set, and so it is important to have efficient algorithms for solving

介绍了三种二次
规划的求解算法

Figure 7.4 Illustration of the ν -SVM applied to a nonseparable data set in two dimensions. The support vectors are indicated by circles.



the quadratic programming problem. We first note that the objective function $\tilde{L}(\mathbf{a})$ given by (7.10) or (7.32) is quadratic and so any local optimum will also be a global optimum provided the constraints define a convex region (which they do as a consequence of being linear). Direct solution of the quadratic programming problem using traditional techniques is often infeasible due to the demanding computation and memory requirements, and so more practical approaches need to be found. The technique of *chunking* (Vapnik, 1982) exploits the fact that the value of the Lagrangian is unchanged if we remove the rows and columns of the kernel matrix corresponding to Lagrange multipliers that have value zero. This allows the full quadratic programming problem to be broken down into a series of smaller ones, whose goal is eventually to identify all of the nonzero Lagrange multipliers and discard the others. Chunking can be implemented using *protected conjugate gradients* (Burges, 1998). Although chunking reduces the size of the matrix in the quadratic function from the number of data points squared to approximately the number of nonzero Lagrange multipliers squared, even this may be too big to fit in memory for large-scale applications. *Decomposition methods* (Osuna *et al.*, 1996) also solve a series of smaller quadratic programming problems but are designed so that each of these is of a fixed size, and so the technique can be applied to arbitrarily large data sets. However, it still involves numerical solution of quadratic programming subproblems and these can be problematic and expensive. One of the most popular approaches to training support vector machines is called *sequential minimal optimization, or SMO* (Platt, 1999). It takes the concept of chunking to the extreme limit and considers just two Lagrange multipliers at a time. In this case, the subproblem can be solved analytically, thereby avoiding numerical quadratic programming altogether. Heuristics are given for choosing the pair of Lagrange multipliers to be considered at each step. In practice, SMO is found to have a scaling with the number of data points that is somewhere between linear and quadratic depending on the particular application.

We have seen that kernel functions correspond to inner products in feature spaces that can have high, or even infinite, dimensionality. By working directly in terms of the kernel function, without introducing the feature space explicitly, (it might therefore seem that) support vector machines somehow manage to avoid the curse of di-

不只是对SVM,其他
在kernel machine 也是如此。它们可以在基向
量为高维甚至无穷维时
进行学习,并不意味着
它们克服了维度诅咒
问题,而是因为高维特征之
间存在约束,有效维度并不
很高。

Section 1.4

事实并非如此

mensionality. (This is not the case) however, because there are constraints amongst the feature values that restrict the effective dimensionality of feature space. To see this consider a simple second-order polynomial kernel that we can expand in terms of its components

$$\begin{aligned}
 k(\mathbf{x}, \mathbf{z}) &= (1 + \mathbf{x}^T \mathbf{z})^2 = (1 + x_1 z_1 + x_2 z_2)^2 \\
 &= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 x_2 z_2 + x_2^2 z_2^2 \\
 &= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\
 &= \phi(\mathbf{x})^T \phi(\mathbf{z}).
 \end{aligned} \tag{7.42}$$

This kernel function therefore represents an inner product in a feature space having six dimensions, in which the mapping from input space to feature space is described by the vector function $\phi(\mathbf{x})$. However, the coefficients weighting these different features are constrained to have specific forms. Thus any set of points in the original two-dimensional space \mathbf{x} would be constrained to lie exactly on a two-dimensional nonlinear manifold embedded in the six-dimensional feature space.

We have already highlighted the fact that the support vector machine does not provide probabilistic outputs but instead makes classification decisions for new input vectors. Veropoulos *et al.* (1999) discuss modifications to the SVM to allow the trade-off between false positive and false negative errors to be controlled. However, if we wish to use the SVM as a module in a larger probabilistic system, then probabilistic predictions of the class label t for new inputs \mathbf{x} are required.

To address this issue, Platt (2000) has proposed fitting a logistic sigmoid to the outputs of a previously trained support vector machine. Specifically, the required conditional probability is assumed to be of the form

$$y(\mathbf{x}) \text{给出了间隔大小, 它显然与间隔相关, 间隔越大, 为某类别概率也越大} \quad p(t=1|\mathbf{x}) = \sigma(Ay(\mathbf{x}) + B) \tag{7.43}$$

但 $y(\mathbf{x})$ 本身却没有概率的含义, 无法通过直接地推导获得概率 $P(t|\mathbf{x})$, 其与概率之间
where $y(\mathbf{x})$ is defined by (7.1). Values for the parameters A and B are found by minimizing the cross-entropy error function defined by a training set consisting of pairs of values $y(\mathbf{x}_n)$ and t_n . The data used to fit the sigmoid needs to be independent of that used to train the original SVM in order to avoid severe over-fitting. This two-stage approach is equivalent to assuming that the output $y(\mathbf{x})$ of the support vector machine represents the log-odds of \mathbf{x} belonging to class $t = 1$. Because the SVM training procedure is not specifically intended to encourage this, the SVM can give a poor approximation to the posterior probabilities (Tipping, 2001).

7.1.2 Relation to logistic regression

As with the separable case, we can re-cast the SVM for nonseparable distributions in terms of the minimization of a regularized error function. This will also allow us to highlight similarities, and differences, compared to the logistic regression model.

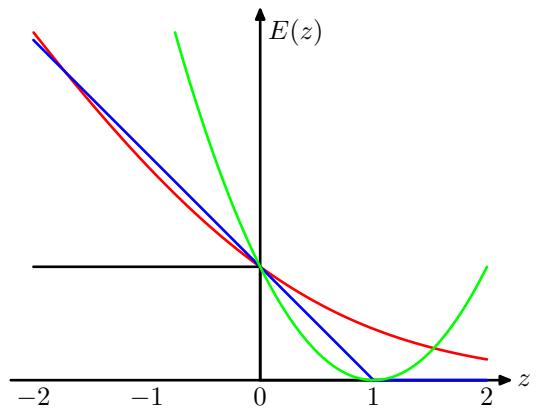
We have seen that for data points that are on the correct side of the margin boundary, and which therefore satisfy $y_n t_n \geq 1$, we have $\xi_n = 0$, and for the

使用SVM构造
概率模型而尝试, 效果并不好。
而映射关系需通过学习得到。同样地, 直接使用
 $p(t=1|\mathbf{x}) = \Gamma(y(\mathbf{x}))$ 也是不合适的, 因为对 $y(\mathbf{x})$ 的学习
并不是通过 $\max \Gamma(y(\mathbf{x}))$
进行的。 $y(\mathbf{x})$ 的学习基于
最大化间隔, 得到的结果是与间隔具有相关性。

Section 4.3.2

性, 但是不是 $\Gamma(y(\mathbf{x}))$
的量化表示未知的
需要引入参数作进一步的拟合。此时, 学习分为 SVM 的学习 + 参数 A, B 的学习 2 个阶段,
而若基于最大化间隔, 直接 $\max \Gamma(y(\mathbf{x}))$ 则是一个 end-to-end 的学习过程。

Figure 7.5 Plot of the ‘hinge’ error function used in support vector machines, shown in blue, along with the error function for logistic regression, rescaled by a factor of $1/\ln(2)$ so that it passes through the point $(0, 1)$, shown in red. Also shown are the misclassification error in black and the squared error in green.



SUM 与 损失 正则化 项 的 形式

remaining points we have $\xi_n = 1 - y_n t_n$. [Thus the objective function (7.21) can be written (up to an overall multiplicative constant) in the form

$$\sum_{n=1}^N E_{\text{SV}}(y_n t_n) + \lambda \|\mathbf{w}\|^2 \quad (7.44)$$

where $\lambda = (2C)^{-1}$, and $E_{\text{SV}}(\cdot)$ is the *hinge* error function defined by

$$E_{\text{SV}}(y_n t_n) = [1 - y_n t_n]_+ \quad (7.45)$$

where $[\cdot]_+$ denotes the positive part. The hinge error function, so-called because of its shape, is plotted in Figure 7.5. It can be viewed as an approximation to the misclassification error, i.e., the error function that ideally we would like to minimize, which is also shown in Figure 7.5.]

logistic regression model [When we considered the logistic regression model in Section 4.3.2, we found it convenient to work with target variable $t \in \{0, 1\}$. For comparison with the support vector machine, we first reformulate maximum likelihood logistic regression using the target variable $t \in \{-1, 1\}$. To do this, we note that $p(t = 1|y) = \sigma(y)$ where $y(\mathbf{x})$ is given by (7.1), and $\sigma(y)$ is the logistic sigmoid function defined by (4.59). It follows that $p(t = -1|y) = 1 - \sigma(y) = \sigma(-y)$, where we have used the properties of the logistic sigmoid function, and so we can write

$$p(t|y) = \sigma(yt). \quad (7.46)$$

Exercise 7.6 From this we can construct an error function by taking the negative logarithm of the likelihood function that, with a quadratic regularizer, takes the form

$$\sum_{n=1}^N E_{\text{LR}}(y_n t_n) + \lambda \|\mathbf{w}\|^2. \quad (7.47)$$

where

$$E_{\text{LR}}(yt) = \ln(1 + \exp(-yt)). \quad (7.48)$$

For comparison with other error functions, we can divide by $\ln(2)$ so that the error function passes through the point $(0, 1)$. This rescaled error function is also plotted in Figure 7.5 and we see that it has a similar form to the support vector error function. The key difference is that the flat region in $E_{SV}(yt)$ leads to sparse solutions 3]

Both the logistic error and the hinge loss can be viewed as continuous approximations to the misclassification error. Another continuous error function that has sometimes been used to solve classification problems is the squared error, which is again plotted in Figure 7.5. It has the property, however, of placing increasing emphasis on data points that are correctly classified but that are a long way from the decision boundary on the correct side. Such points will be strongly weighted at the expense of misclassified points, and so if the objective is to minimize the misclassification rate, then a monotonically decreasing error function would be a better choice.

7.1.3 Multiclass SVMs

The support vector machine is fundamentally a two-class classifier. In practice, however, we often have to tackle problems involving $K > 2$ classes. Various methods have therefore been proposed for combining multiple two-class SVMs in order to build a multiclass classifier.

One commonly used approach (Vapnik, 1998) is to construct K separate SVMs, in which the k^{th} model $y_k(\mathbf{x})$ is trained using the data from class C_k as the positive examples and the data from the remaining $K - 1$ classes as the negative examples. This is known as the one-versus-the-rest approach. However, in Figure 4.2 we saw that using the decisions of the individual classifiers can lead to inconsistent results in which an input is assigned to multiple classes simultaneously. This problem is sometimes addressed by making predictions for new inputs \mathbf{x} using

$$y(\mathbf{x}) = \max_k y_k(\mathbf{x}). \quad (7.49)$$

Unfortunately, this heuristic approach suffers from the problem that the different classifiers were trained on different tasks, and there is no guarantee that the real-valued quantities $y_k(\mathbf{x})$ for different classifiers will have appropriate scales.

Another problem with the one-versus-the-rest approach is that the training sets are imbalanced. For instance, if we have ten classes each with equal numbers of training data points, then the individual classifiers are trained on data sets comprising 90% negative examples and only 10% positive examples, and the symmetry of the original problem is lost. A variant of the one-versus-the-rest scheme was proposed by Lee *et al.* (2001) who modify the target values so that the positive class has target +1 and the negative class has target $-1/(K - 1)$. 3

Weston and Watkins (1999) define a single objective function for training all K SVMs simultaneously, based on maximizing the margin from each to remaining classes. However, this can result in much slower training because, instead of solving K separate optimization problems each over N data points with an overall cost of $O(KN^2)$, a single optimization problem of size $(K - 1)N$ must be solved giving an overall cost of $O(K^2N^2)$. 3

③ Another approach is to train $K(K - 1)/2$ different 2-class SVMs on all possible pairs of classes, and then to classify test points according to which class has the highest number of ‘votes’, an approach that is sometimes called *one-versus-one*. Again, we saw in Figure 4.2 that this can lead to ambiguities in the resulting classification. Also, for large K this approach requires significantly more training time than the one-versus-the-rest approach. Similarly, to evaluate test points, significantly more computation is required.

The latter problem can be alleviated by organizing the pairwise classifiers into a directed acyclic graph (not to be confused with a probabilistic graphical model) leading to the DAGSVM (Platt *et al.*, 2000). For K classes, the DAGSVM has a total of $K(K - 1)/2$ classifiers, and to classify a new test point only $K - 1$ pairwise classifiers need to be evaluated, with the particular classifiers used depending on which path through the graph is traversed. ③

④ A different approach to multiclass classification, based on error-correcting output codes, was developed by Dietterich and Bakiri (1995) and applied to support vector machines by Allwein *et al.* (2000). This can be viewed as a generalization of the voting scheme of the one-versus-one approach in which more general partitions of the classes are used to train the individual classifiers. The K classes themselves are represented as particular sets of responses from the two-class classifiers chosen, and together with a suitable decoding scheme, this gives robustness to errors and to ambiguity in the outputs of the individual classifiers. ③ Although the application of SVMs to multiclass classification problems remains an open issue, in practice the one-versus-the-rest approach is the most widely used in spite of its ad-hoc formulation and its practical limitations.

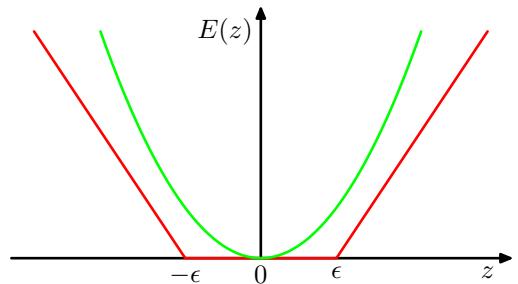
支持向量机在概率密度估计这一大监督学习中的应用

There are also *single-class* support vector machines, which solve an unsupervised learning problem related to probability density estimation. Instead of modelling the density of data, however, these methods aim to find a smooth boundary enclosing a region of high density. The boundary is chosen to represent a quantile of the density, that is, the probability that a data point drawn from the distribution will land inside that region is given by a fixed number between 0 and 1 that is specified in advance. This is a more restricted problem than estimating the full density but may be sufficient in specific applications. Two approaches to this problem using support vector machines have been proposed. The algorithm of Schölkopf *et al.* (2001) tries to find a hyperplane that separates all but a fixed fraction ν of the training data from the origin while at the same time maximizing the distance (margin) of the hyperplane from the origin, while Tax and Duin (1999) look for the smallest sphere in feature space that contains all but a fraction ν of the data points. For kernels $k(\mathbf{x}, \mathbf{x}')$ that are functions only of $\mathbf{x} - \mathbf{x}'$, the two algorithms are equivalent. ③

7.1.4 SVMs for regression

We now extend support vector machines to regression problems while at the same time preserving the property of sparseness. In simple linear regression, we

Figure 7.6 Plot of an ϵ -insensitive error function (in red) in which the error increases linearly with distance beyond the insensitive region. Also shown for comparison is the quadratic error function (in green).



minimize a regularized error function given by

$$\frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2. \quad (7.50)$$

To obtain sparse solutions, the quadratic error function is replaced by an *ϵ -insensitive error function* (Vapnik, 1995), which gives zero error if the absolute difference between the prediction $y(\mathbf{x})$ and the target t is less than ϵ where $\epsilon > 0$. A simple example of an ϵ -insensitive error function, having a linear cost associated with errors outside the insensitive region, is given by

$$E_\epsilon(y(\mathbf{x}) - t) = \begin{cases} 0, & \text{if } |y(\mathbf{x}) - t| < \epsilon; \\ |y(\mathbf{x}) - t| - \epsilon, & \text{otherwise} \end{cases} \quad (7.51)$$

and is illustrated in Figure 7.6.

We therefore minimize a regularized error function given by

$$C \sum_{n=1}^N E_\epsilon(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.52)$$

where $y(\mathbf{x})$ is given by (7.1). By convention the (inverse) regularization parameter, denoted C , appears in front of the error term.

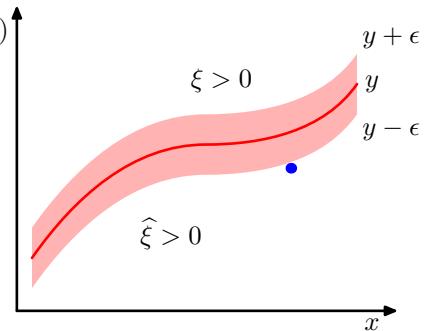
As before, we can re-express the optimization problem by introducing slack variables. For each data point \mathbf{x}_n , we now need two slack variables $\xi_n \geq 0$ and $\hat{\xi}_n \geq 0$, where $\xi_n > 0$ corresponds to a point for which $t_n > y(\mathbf{x}_n) + \epsilon$, and $\hat{\xi}_n > 0$ corresponds to a point for which $t_n < y(\mathbf{x}_n) - \epsilon$, as illustrated in Figure 7.7.

The condition for a target point to lie inside the ϵ -tube is that $y_n - \epsilon \leq t_n \leq y_n + \epsilon$, where $y_n = y(\mathbf{x}_n)$. Introducing the slack variables allows points to lie outside the tube provided the slack variables are nonzero, and the corresponding conditions are

$$t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n \quad (7.53)$$

$$t_n \geq y(\mathbf{x}_n) - \epsilon - \hat{\xi}_n. \quad (7.54)$$

Figure 7.7 Illustration of SVM regression, showing the regression curve together with the ϵ -insensitive ‘tube’. Also shown are examples of the slack variables ξ and $\hat{\xi}$. Points above the ϵ -tube have $\xi > 0$ and $\hat{\xi} = 0$, points below the ϵ -tube have $\xi = 0$ and $\hat{\xi} > 0$, and points inside the ϵ -tube have $\xi = \hat{\xi} = 0$.



The error function for support vector regression can then be written as

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.55)$$

which must be minimized subject to the constraints $\xi_n \geq 0$ and $\hat{\xi}_n \geq 0$ as well as (7.53) and (7.54). This can be achieved by introducing Lagrange multipliers $a_n \geq 0$, $\hat{a}_n \geq 0$, $\mu_n \geq 0$, and $\hat{\mu}_n \geq 0$ and optimizing the Lagrangian

$$\begin{aligned} L = & C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \\ & - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n). \end{aligned} \quad (7.56)$$

We now substitute for $y(\mathbf{x})$ using (7.1) and then set the derivatives of the Lagrangian with respect to \mathbf{w} , b , ξ_n , and $\hat{\xi}_n$ to zero, giving

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N (a_n - \hat{a}_n) \phi(\mathbf{x}_n) \quad (7.57)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0 \quad (7.58)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n + \mu_n = C \quad (7.59)$$

$$\frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{a}_n + \hat{\mu}_n = C. \quad (7.60)$$

Using these results to eliminate the corresponding variables from the Lagrangian, we see that the dual problem involves maximizing

$$\begin{aligned}\tilde{L}(\mathbf{a}, \widehat{\mathbf{a}}) &= -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \widehat{a}_n)(a_m - \widehat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \\ &\quad - \epsilon \sum_{n=1}^N (a_n + \widehat{a}_n) + \sum_{n=1}^N (a_n - \widehat{a}_n)t_n\end{aligned}\quad (7.61)$$

with respect to $\{a_n\}$ and $\{\widehat{a}_n\}$, where we have introduced the kernel $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$. Again, this is a constrained maximization, and to find the constraints we note that $a_n \geq 0$ and $\widehat{a}_n \geq 0$ are both required because these are Lagrange multipliers. Also $\mu_n \geq 0$ and $\widehat{\mu}_n \geq 0$ together with (7.59) and (7.60), require $a_n \leq C$ and $\widehat{a}_n \leq C$, and so again we have the box constraints

$$0 \leq a_n \leq C \quad (7.62)$$

$$0 \leq \widehat{a}_n \leq C \quad (7.63)$$

together with the condition (7.58).

Substituting (7.57) into (7.1), we see that predictions for new inputs can be made using

$$y(\mathbf{x}) = \sum_{n=1}^N (a_n - \widehat{a}_n) k(\mathbf{x}, \mathbf{x}_n) + b \quad (7.64)$$

which is again expressed in terms of the kernel function.

The corresponding Karush-Kuhn-Tucker (KKT) conditions, which state that at the solution the product of the dual variables and the constraints must vanish, are given by

$$a_n(\epsilon + \xi_n + y_n - t_n) = 0 \quad (7.65)$$

$$\widehat{a}_n(\epsilon + \widehat{\xi}_n - y_n + t_n) = 0 \quad (7.66)$$

$$(C - a_n)\xi_n = 0 \quad (7.67)$$

$$(C - \widehat{a}_n)\widehat{\xi}_n = 0. \quad (7.68)$$

From these we can obtain several useful results. First of all, we note that a coefficient a_n can only be nonzero if $\epsilon + \xi_n + y_n - t_n = 0$, which implies that the data point either lies on the upper boundary of the ϵ -tube ($\xi_n = 0$) or lies above the upper boundary ($\xi_n > 0$). Similarly, a nonzero value for \widehat{a}_n implies $\epsilon + \widehat{\xi}_n - y_n + t_n = 0$, and such points must lie either on or below the lower boundary of the ϵ -tube.

Furthermore, the two constraints $\epsilon + \xi_n + y_n - t_n = 0$ and $\epsilon + \widehat{\xi}_n - y_n + t_n = 0$ are incompatible, as is easily seen by adding them together and noting that ξ_n and $\widehat{\xi}_n$ are nonnegative while ϵ is strictly positive, and so for every data point \mathbf{x}_n , either a_n or \widehat{a}_n (or both) must be zero.

The support vectors are those data points that contribute to predictions given by (7.64), in other words those for which either $a_n \neq 0$ or $\widehat{a}_n \neq 0$. These are points that lie on the boundary of the ϵ -tube or outside the tube. All points within the tube have

$a_n = \hat{a}_n = 0$. We again have a sparse solution, and the only terms that have to be evaluated in the predictive model (7.64) are those that involve the support vectors.

The parameter b can be found by considering a data point for which $0 < a_n < C$, which from (7.67) must have $\xi_n = 0$, and from (7.65) must therefore satisfy $\epsilon + y_n - t_n = 0$. Using (7.1) and solving for b , we obtain

$$\begin{aligned} b &= t_n - \epsilon - \mathbf{w}^T \phi(\mathbf{x}_n) \\ &= t_n - \epsilon - \sum_{m=1}^N (a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \end{aligned} \quad (7.69)$$

where we have used (7.57). We can obtain an analogous result by considering a point for which $0 < \hat{a}_n < C$. In practice, it is better to average over all such estimates of b .

νSVM for regression [As with the classification case, there is an alternative formulation of the SVM for regression in which the parameter governing complexity has a more intuitive interpretation (Schölkopf *et al.*, 2000). In particular, instead of fixing the width ϵ of the insensitive region, we fix instead a parameter ν that bounds the fraction of points lying outside the tube. This involves maximizing

$$\begin{aligned} \tilde{L}(\mathbf{a}, \hat{\mathbf{a}}) &= -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \\ &\quad + \sum_{n=1}^N (a_n - \hat{a}_n)t_n \end{aligned} \quad (7.70)$$

subject to the constraints

$$0 \leq a_n \leq C/N \quad (7.71)$$

$$0 \leq \hat{a}_n \leq C/N \quad (7.72)$$

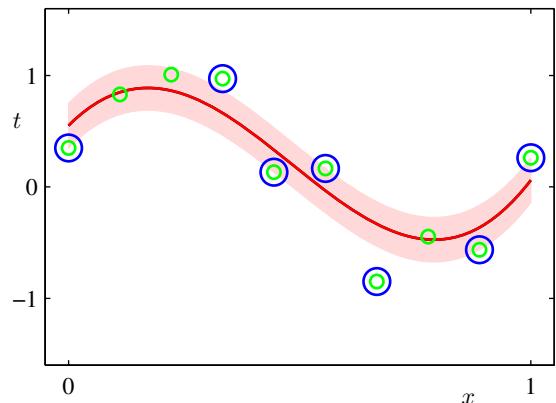
$$\sum_{n=1}^N (a_n - \hat{a}_n) = 0 \quad (7.73)$$

$$\sum_{n=1}^N (a_n + \hat{a}_n) \leq \nu C. \quad (7.74)$$

It can be shown that there are at most νN data points falling outside the insensitive tube, while at least νN data points are support vectors and so lie either on the tube or outside it.]

The use of a support vector machine to solve a regression problem is illustrated using the sinusoidal data set in Figure 7.8. Here the parameters ν and C have been chosen by hand. In practice, their values would typically be determined by cross-validation.

Figure 7.8 Illustration of the ν -SVM for regression applied to the sinusoidal synthetic data set using Gaussian kernels. The predicted regression curve is shown by the red line, and the ϵ -insensitive tube corresponds to the shaded region. Also, the data points are shown in green, and those with support vectors are indicated by blue circles.



7.1.5 Computational learning theory

只简单地介绍了计算学习理论
并未阐述其与SVM的具体联系

Historically, support vector machines have largely been motivated and analysed using a theoretical framework known as *computational learning theory*, also sometimes called *statistical learning theory* (Anthony and Biggs, 1992; Kearns and Vazirani, 1994; Vapnik, 1995; Vapnik, 1998). This has its origins with Valiant (1984) who formulated the *probably approximately correct*, or PAC, learning framework. The goal of the PAC framework is to understand how large a data set needs to be in order to give good generalization. It also gives bounds for the computational cost of learning, although we do not consider these here.

Suppose that a data set \mathcal{D} of size N is drawn from some joint distribution $p(\mathbf{x}, \mathbf{t})$ where \mathbf{x} is the input variable and \mathbf{t} represents the class label, and that we restrict attention to ‘noise free’ situations in which the class labels are determined by some (unknown) deterministic function $\mathbf{t} = \mathbf{g}(\mathbf{x})$. In PAC learning we say that a function $\mathbf{f}(\mathbf{x}; \mathcal{D})$, drawn from a space \mathcal{F} of such functions on the basis of the training set \mathcal{D} , has good generalization if its expected error rate is below some pre-specified threshold ϵ , so that

$$\text{是范斯基的 } \mathbf{x}, \mathbf{t} \quad \text{是训练样本的训练集} \quad \text{而非测试集} \quad \mathbb{E}_{\mathbf{x}, \mathbf{t}} [I(\mathbf{f}(\mathbf{x}; \mathcal{D}) \neq \mathbf{t})] < \epsilon \quad (7.75)$$

中括号 \mathbf{x}, \mathbf{t}

where $I(\cdot)$ is the indicator function, and the expectation is with respect to the distribution $p(\mathbf{x}, \mathbf{t})$. The quantity on the left-hand side is a random variable, because it depends on the training set \mathcal{D} , and the PAC framework requires that (7.75) holds, with probability greater than $1 - \delta$, for a data set \mathcal{D} drawn randomly from $p(\mathbf{x}, \mathbf{t})$. Here δ is another pre-specified parameter, and the terminology ‘probably approximately correct’ comes from the requirement that with high probability (greater than $1 - \delta$), the error rate be small (less than ϵ). For a given choice of model space \mathcal{F} , and for given parameters ϵ and δ , PAC learning aims to provide bounds on the minimum size N of data set needed to meet this criterion. A key quantity in PAC learning is the *Vapnik-Chervonenkis dimension*, or VC dimension, which provides a measure of the complexity of a space of functions, and which allows the PAC framework to be extended to spaces containing an infinite number of functions.

The bounds derived within the PAC framework are often described as worst-

case, because they apply to *any* choice for the distribution $p(\mathbf{x}, \mathbf{t})$, so long as both the training and the test examples are drawn (independently) from the same distribution, and for *any* choice for the function $f(\mathbf{x})$ so long as it belongs to \mathcal{F} . In real-world applications of machine learning, we deal with distributions that have significant regularity, for example in which large regions of input space carry the same class label. As a consequence of the lack of any assumptions about the form of the distribution, the PAC bounds are very conservative, in other words they strongly over-estimate the size of data sets required to achieve a given generalization performance. For this reason, PAC bounds have found few, if any, practical applications.

One attempt to improve the tightness of the PAC bounds is the *PAC-Bayesian* framework (McAllester, 2003), which considers a distribution over the space \mathcal{F} of functions, somewhat analogous to the prior in a Bayesian treatment. This still considers any possible choice for $p(\mathbf{x}, \mathbf{t})$, and so although the bounds are tighter, they are still very conservative.

7.2. Relevance Vector Machines

SVM缺点 Support vector machines have been used in a variety of classification and regression applications. Nevertheless, they suffer from a number of limitations, several of which have been highlighted already in this chapter. In particular, the outputs of an SVM represent decisions rather than posterior probabilities. Also, the SVM was originally formulated for two classes, and the extension to $K > 2$ classes is problematic. There is a complexity parameter C , or ν (as well as a parameter ϵ in the case of regression), that must be found using a hold-out method such as cross-validation. Finally, predictions are expressed as linear combinations of kernel functions that are centred on training data points and that are required to be positive definite. **核函数保证**

RVM The *relevance vector machine* or RVM (Tipping, 2001) is a Bayesian sparse kernel technique for regression and classification that shares many of the characteristics of the SVM whilst avoiding its principal limitations. Additionally, it typically leads to much sparser models resulting in correspondingly faster performance on test data whilst maintaining comparable generalization error. **严格, 需要使用核函数**

In contrast to the SVM we shall find it more convenient to introduce the regression form of the RVM first and then consider the extension to classification tasks. **卫足。另外**

7.2.1 RVM for regression

模型介绍 [The relevance vector machine for regression is a linear model of the form studied in Chapter 3 but with a modified prior that results in sparse solutions. The model defines a conditional distribution for a real-valued target variable t , given an input vector \mathbf{x} , which takes the form

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}), \beta^{-1}) \quad (7.76)$$

where $\beta = \sigma^{-2}$ is the noise precision (inverse noise variance), and the mean is given by a linear model of the form

$$y(\mathbf{x}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (7.77)$$

with fixed nonlinear basis functions $\phi_i(\mathbf{x})$, which will typically include a constant term so that the corresponding weight parameter represents a ‘bias’.

The relevance vector machine is a specific instance of this model, which is intended to mirror the structure of the support vector machine. In particular, the basis functions are given by kernels, with one kernel associated with each of the data points from the training set. The general expression (7.77) then takes the SVM-like form

RVM模型式(7.78)是式(7.77)
的特例,下文叫惟导基式
(7.77),具有-般性,因此对式(7.78)成立

$$y(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}, \mathbf{x}_n) + b \quad \begin{array}{l} \text{注意于SUM中不同的是 } x_n \\ \text{并不只是训练数据} \\ \text{而与下文不同不是训练数据} \end{array} \quad (7.78)$$

where b is a bias parameter. The number of parameters in this case is $M = N + 1$, and $y(\mathbf{x})$ has the same form as the predictive model (7.64) for the SVM, except that the coefficients a_n are here denoted w_n . It should be emphasized that the subsequent analysis is valid for arbitrary choices of basis function, and for generality we shall work with the form (7.77). In contrast to the SVM, there is no restriction to positive-definite kernels, nor are the basis functions tied in either number or location to the training data points.

// Suppose we are given a set of N observations of the input vector \mathbf{x} , which we denote collectively by a data matrix \mathbf{X} whose n^{th} row is \mathbf{x}_n^T with $n = 1, \dots, N$. The corresponding target values are given by $\mathbf{t} = (t_1, \dots, t_N)^T$. Thus, the likelihood function is given by

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta) \quad (7.79)$$

Next we introduce a prior distribution over the parameter vector \mathbf{w} and as in Chapter 3, we shall consider a zero-mean Gaussian prior. However, the key difference in the RVM is that we introduce a separate hyperparameter α_i for each of the weight parameters w_i instead of a single shared hyperparameter. Thus the weight prior takes the form

由下文可知,参考稀疏解
关键在于采用了式(7.80)
形式优先,遍优化

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha_i^{-1}) \quad \begin{array}{l} \text{参数 } \alpha, \beta, \text{是通过得稀疏} \\ \text{而建立} \end{array} \quad (7.80)$$

where α_i represents the precision of the corresponding parameter w_i , and $\boldsymbol{\alpha}$ denotes $(\alpha_1, \dots, \alpha_M)^T$. We shall see that, when we maximize the evidence with respect to these hyperparameters, a significant proportion of them go to infinity, and the corresponding weight parameters have posterior distributions that are concentrated at zero. The basis functions associated with these parameters therefore play no role

式(7.80)的 auto
relevance determination
(ARD) prior

in the predictions made by the model and so are effectively pruned out, resulting in a sparse model.

模型参数 [Using the result (3.49) for linear regression models, we see that the posterior **模型证据** distribution for the weights is again Gaussian and takes the form

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \boldsymbol{\alpha}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \Sigma) \quad (7.81)$$

where the mean and covariance are given by

$$\mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t} \quad (7.82)$$

$$\boldsymbol{\Sigma} = (\mathbf{A} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \quad \begin{array}{l} \text{for } i=1, \dots, N \text{ and} \\ \mathbf{z}_{ii} = 1 \text{ for } n=1, \dots, N \end{array} \quad (7.83)$$

where $\boldsymbol{\Phi}$ is the $N \times M$ design matrix with elements $\Phi_{ni} = \phi_i(\mathbf{x}_n)$, and $\mathbf{A} = \text{diag}(\alpha_i)$. Note that in the specific case of the model (7.78), we have $\boldsymbol{\Phi} = \mathbf{K}$, where \mathbf{K} is the symmetric $(N+1) \times (N+1)$ kernel matrix with elements $k(\mathbf{x}_n, \mathbf{x}_m)$.

The values of $\boldsymbol{\alpha}$ and β are determined using type-2 maximum likelihood, also known as the *evidence approximation*, in which we maximize the marginal likelihood function obtained by integrating out the weight parameters

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\boldsymbol{\alpha}) d\mathbf{w}. \quad (7.84)$$

Section 3.5

Exercise 7.10

Because this represents the convolution of two Gaussians, it is readily evaluated to give the log marginal likelihood in the form

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) &= \ln \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \\ &= -\frac{1}{2} \{ N \ln(2\pi) + \ln |\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} \} \end{aligned} \quad (7.85)$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$, and we have defined the $N \times N$ matrix \mathbf{C} given by

$$\mathbf{C} = \beta^{-1} \mathbf{I} + \begin{pmatrix} \text{diag}(\alpha_i) \\ \mathbf{\Phi} \mathbf{A}^{-1} \mathbf{\Phi}^T \end{pmatrix} \quad \begin{array}{l} \text{就是相应矩阵} \\ \text{核函数的Gram matrix,} \end{array} \quad (7.86)$$

Exercise 7.12

Our goal is now to maximize (7.85) with respect to the hyperparameters $\boldsymbol{\alpha}$ and β . This requires only a small modification to the results obtained in Section 3.5 for the evidence approximation in the linear regression model. Again, we can identify two approaches. In the first, we simply set the required derivatives of the marginal likelihood to zero and obtain the following re-estimation equations

$$\alpha_i^{\text{new}} = \frac{\gamma_i}{m_i^2} \quad \begin{array}{l} \text{只立这里并采} \\ \text{用} \end{array} \quad (7.87)$$

$$(\beta^{\text{new}})^{-1} = \frac{\|\mathbf{t} - \boldsymbol{\Phi} \mathbf{m}\|^2}{N - \sum_i \gamma_i} \quad \begin{array}{l} \text{介绍, 实际上与相} \\ \text{反列是新列,} \end{array} \quad (7.88)$$

where m_i is the i^{th} component of the posterior mean \mathbf{m} defined by (7.82). The quantity γ_i measures how well the corresponding parameter w_i is determined by the data and is defined by

Section 3.5.3

即底层逻辑是一样的
只是表现形式不同。

$$\gamma_i = 1 - \alpha_i \Sigma_{ii} \quad (7.89)$$

in which Σ_{ii} is the i^{th} diagonal component of the posterior covariance Σ given by (7.83). Learning therefore proceeds by choosing initial values for α and β , evaluating the mean and covariance of the posterior using (7.82) and (7.83), respectively, and then alternately re-estimating the hyperparameters, using (7.87) and (7.88), and re-estimating the posterior mean and covariance, using (7.82) and (7.83), until a suitable convergence criterion is satisfied.

// The second approach is to use the EM algorithm, and is discussed in Section 9.3.4. These two approaches to finding the values of the hyperparameters that maximize the evidence are formally equivalent. Numerically, however, it is found that the direct optimization approach corresponding to (7.87) and (7.88) gives somewhat faster convergence (Tipping, 2001). }

Exercise 9.23

Section 7.2.2

As a result of the optimization, we find that a proportion of the hyperparameters $\{\alpha_i\}$ are driven to large (in principle infinite) values, and so the weight parameters w_i corresponding to these hyperparameters have posterior distributions with mean and variance both zero. Thus those parameters, and the corresponding basis functions $\phi_i(x)$, are removed from the model and play no role in making predictions for new inputs. In the case of models of the form (7.78), the inputs x_n corresponding to the remaining nonzero weights are called *relevance vectors*, because they are identified through the mechanism of automatic relevance determination, and are analogous to the support vectors of an SVM. It is worth emphasizing, however, that this mechanism for achieving sparsity in probabilistic models through automatic relevance determination is quite general and can be applied to any model expressed as an adaptive linear combination of basis functions.]

Exercise 7.14

预测 [Having found values α^* and β^* for the hyperparameters that maximize the marginal likelihood, we can evaluate the predictive distribution over t for a new input x . Using (7.76) and (7.81), this is given by

$$\begin{aligned} p(t|x, \mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) &= \int p(t|\mathbf{w}, \beta^*) p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}^T \phi(x), \sigma^2(x)). \end{aligned} \quad (7.90)$$

Thus the predictive mean is given by (7.76) with w set equal to the posterior mean m , and the variance of the predictive distribution is given by

$$\sigma^2(x) = (\beta^*)^{-1} + \phi(x)^T \Sigma \phi(x) \quad (7.91)$$

where Σ is given by (7.83) in which α and β are set to their optimized values α^* and β^* . This is just the familiar result (3.59) obtained in the context of linear regression.]

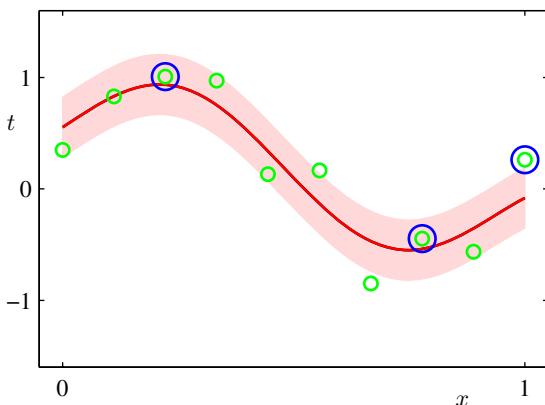
RIMM译 [Recall that for localized basis functions, the predictive variance for linear regression models becomes small in regions of input space where there are no basis functions. In the case of an RVM with the basis functions centred on data points, the model will therefore become increasingly certain of its predictions when extrapolating outside the domain of the data (Rasmussen and Quiñonero-Candela, 2005), which of course is undesirable. The predictive distribution in Gaussian process regression does not

个人认为这是说法有误：前面说若RVM形式(7.78)中的基函数使用
localized basis function，外插时结果会显得不合理。同样地，对高斯过程
而言，若其核函数对应的基函数也是localized，显然也会遇到这种问题。
因为高斯过程只是表现形式不同，它们的本质属性是一样的。这里之所

Section 6.4.2

Figure 7.9
-SVM 与对比时, RVM
即式(7.18)中的核函数
应同 SVM, x_n 也是训练
数据, 即对相同的数据
使用不同的方法进行学习
(SVM or RVM), 否则若
模型都不同, 则训练结
果就没有比较的意义了。

Illustration of RVM regression using the same data set, and the same Gaussian kernel functions, as used in Figure 7.8 for the ν -SVM regression model. The mean of the predictive distribution for the RVM is shown by the red line, and the one standard-deviation predictive distribution is shown by the shaded region. Also, the data points are shown in green, and the relevance vectors are indicated by blue circles. Note that there are only 3 relevance vectors compared to 7 support vectors for the ν -SVM in Figure 7.8.



suffer from this problem. However, the computational cost of making predictions with a Gaussian process is typically much higher than with an RVM.

Figure 7.9 shows an example of the RVM applied to the sinusoidal regression data set. Here the noise precision parameter β is also determined through evidence maximization. We see that the number of relevance vectors in the RVM is significantly smaller than the number of support vectors used by the SVM. For a wide range of regression and classification tasks, the RVM is found to give models that are typically an order of magnitude more compact than the corresponding support vector machine, resulting in a significant improvement in the speed of processing on test data. Remarkably, this greater sparsity is achieved with little or no reduction in generalization error compared with the corresponding SVM.

The principal disadvantage of the RVM compared to the SVM is that training involves optimizing a nonconvex function, and training times can be longer than for a comparable SVM. For a model with M basis functions, the RVM requires inversion of a matrix of size $M \times M$, which in general requires $O(M^3)$ computation. In the specific case of the SVM-like model (7.78), we have $M = N + 1$. As we have noted, there are techniques for training SVMs whose cost is roughly quadratic in N . Of course, in the case of the RVM we always have the option of starting with a smaller number of basis functions than $N + 1$. More significantly, in the relevance vector machine the parameters governing complexity and noise variance are determined automatically from a single training run, whereas in the support vector machine the parameters C and ϵ (or ν) are generally found using cross-validation, which involves multiple training runs. Furthermore, in the next section we shall derive an alternative procedure for training the relevance vector machine that improves training speed significantly.

RVM 相比 SVM
更稀疏, 迭代
效果也不差

RVM 与 SVM 训练
时间相近

7.2.2 Analysis of sparsity

We have noted earlier that the mechanism of automatic relevance determination causes a subset of parameters to be driven to zero. We now examine in more detail

给出产生稀疏解的原理，并由此得到起参数优化的更有效的方法

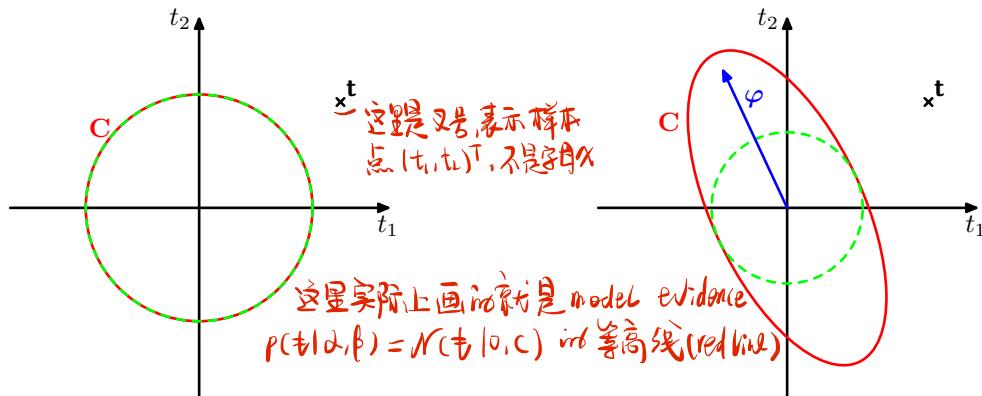


Figure 7.10 Illustration of the mechanism for sparsity in a Bayesian linear regression model, showing a training set vector of target values given by $\mathbf{t} = (t_1, t_2)^T$, indicated by the cross, for a model with one basis vector $\varphi = (\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))^T$, which is poorly aligned with the target data vector \mathbf{t} . On the left we see a model having only isotropic noise, so that $\mathbf{C} = \beta^{-1}\mathbf{I}$, corresponding to $\alpha = \infty$, with β set to its most probable value. On the right we see the same model but with a finite value of α . In each case the red ellipse corresponds to unit Mahalanobis distance, with $|\mathbf{C}|$ taking the same value for both plots, while the dashed green circle shows the contribution arising from the noise term β^{-1} . We see that any finite value of α reduces the probability of the observed data, and so for the most probable solution the basis vector is removed.

contribution

the mechanism of sparsity in the context of the relevance vector machine. In the process, we will arrive at a significantly faster procedure for optimizing the hyperparameters compared to the direct techniques given above.

对称属性
冗余分析

Before proceeding with a mathematical analysis, we first give some informal insight into the origin of sparsity in Bayesian linear models. Consider a data set comprising $N = 2$ observations t_1 and t_2 , together with a model having a single basis function $\phi(\mathbf{x})$, with hyperparameter α , along with isotropic noise having precision β . From (7.85), the marginal likelihood is given by $p(\mathbf{t}|\alpha, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C})$ in which the covariance matrix takes the form

$$\mathbf{C} = \frac{1}{\beta} \mathbf{I} + \frac{1}{\alpha} \varphi \varphi^T \quad (7.92)$$

where φ denotes the N -dimensional vector $(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))^T$, and similarly $\mathbf{t} = (t_1, t_2)^T$. Notice that this is just a zero-mean Gaussian process model over \mathbf{t} with covariance \mathbf{C} . Given a particular observation for \mathbf{t} , our goal is to find α^* and β^* by maximizing the marginal likelihood. We see from Figure 7.10 that, if there is a poor alignment between the direction of φ and that of the training data vector \mathbf{t} , then the corresponding hyperparameter α will be driven to ∞ , and the basis vector will be pruned from the model. This arises because any finite value for α will always assign a lower probability to the data, thereby decreasing the value of the density at \mathbf{t} , provided that β is set to its optimal value. We see that any finite value for α would cause the distribution to be elongated in a direction away from the data, thereby increasing the probability mass in regions away from the observed data and hence reducing the value of the density at the target data vector itself. For the more general case of M

这三处
对应
为

注意，这些基向量不一定是正交的

basis vectors $\varphi_1, \dots, \varphi_M$ a similar intuition holds, namely that if a particular basis vector is poorly aligned with the data vector t , then it is likely to be pruned from the model.

// We now investigate the mechanism for sparsity from a more mathematical perspective, for a general case involving M basis functions. To motivate this analysis we first note that, in the result (7.87) for re-estimating the parameter α_i , the terms on the right-hand side are themselves also functions of α_i . These results therefore represent implicit solutions, and iteration would be required even to determine a single α_i with all other α_j for $j \neq i$ fixed.

This suggests a different approach to solving the optimization problem for the RVM, in which we make explicit all of the dependence of the marginal likelihood (7.85) on a particular α_i and then determine its stationary points explicitly (Faul and Tipping, 2002; Tipping and Faul, 2003). To do this, we first pull out the contribution from α_i in the matrix C defined by (7.86) to give

$$\begin{aligned} C &= \beta^{-1}I + \sum_{j \neq i} \alpha_j^{-1} \varphi_j \varphi_j^T + \alpha_i^{-1} \varphi_i \varphi_i^T \\ &= C_{-i} + \alpha_i^{-1} \varphi_i \varphi_i^T \end{aligned} \quad (7.93)$$

where φ_i denotes the i^{th} column of Φ , in other words the N -dimensional vector with elements $(\phi_i(x_1), \dots, \phi_i(x_N))$, in contrast to ϕ_n , which denotes the n^{th} row of Φ . The matrix C_{-i} represents the matrix C with the contribution from basis function i removed. Using the matrix identities (C.7) and (C.15), the determinant and inverse of C can then be written

$$|C| = |C_{-i}| \underbrace{1 + \alpha_i^{-1} \varphi_i^T C_{-i}^{-1} \varphi_i}_{\text{括号}} \quad (7.94)$$

$$C^{-1} = C_{-i}^{-1} - \frac{C_{-i}^{-1} \varphi_i \varphi_i^T C_{-i}^{-1}}{\alpha_i + \varphi_i^T C_{-i}^{-1} \varphi_i}. \quad (7.95)$$

Exercise 7.15

Using these results, we can then write the log marginal likelihood function (7.85) in the form

$$L(\alpha) = L(\alpha_{-i}) + \lambda(\alpha_i) \quad (7.96)$$

where $L(\alpha_{-i})$ is simply the log marginal likelihood with basis function φ_i omitted, and the quantity $\lambda(\alpha_i)$ is defined by

$$\lambda(\alpha_i) = \frac{1}{2} \left[\ln \alpha_i - \ln (\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right] \quad (7.97)$$

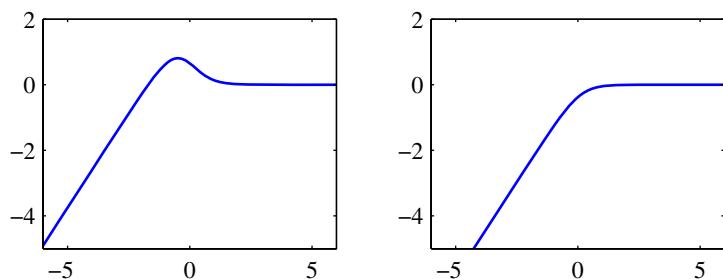
and contains all of the dependence on α_i . Here we have introduced the two quantities

$$s_i = \varphi_i^T C_{-i}^{-1} \varphi_i \quad (7.98)$$

$$q_i = \varphi_i^T C_{-i}^{-1} t. \quad (7.99)$$

Here s_i is called the *sparsity* and q_i is known as the *quality* of φ_i , and as we shall see, a large value of s_i relative to the value of q_i means that the basis function φ_i

Figure 7.11 Plots of the log marginal likelihood $\lambda(\alpha_i)$ versus $\ln \alpha_i$ showing on the left, the single maximum at a finite α_i for $q_i^2 = 4$ and $s_i = 1$ (so that $q_i^2 > s_i$) and on the right, the maximum at $\alpha_i = \infty$ for $q_i^2 = 1$ and $s_i = 2$ (so that $q_i^2 < s_i$).



is more likely to be pruned from the model. The ‘sparsity’ measures the extent to which basis function φ_i overlaps with the other basis vectors in the model, and the ‘quality’ represents a measure of the alignment of the basis vector φ_i with the error between the training set values $\mathbf{t} = (t_1, \dots, t_N)^T$ and the vector \mathbf{y}_{-i} of predictions that would result from the model with the vector φ_i excluded (Tipping and Faul, 2003).}

The stationary points of the marginal likelihood with respect to α_i occur when the derivative

$$\frac{d\lambda(\alpha_i)}{d\alpha_i} = \frac{\alpha_i^{-1}s_i^2 - (q_i^2 - s_i)}{2(\alpha_i + s_i)^2} \quad (7.100)$$

is equal to zero. There are two possible forms for the solution. Recalling that $\alpha_i \geq 0$, we see that if $q_i^2 < s_i$, then $\alpha_i \rightarrow \infty$ provides a solution. Conversely, if $q_i^2 > s_i$, we can solve for α_i to obtain

$$\alpha_i = \frac{s_i^2}{q_i^2 - s_i}. \quad (7.101)$$

These two solutions are illustrated in Figure 7.11. We see that the relative size of the quality and sparsity terms determines whether a particular basis vector will be pruned from the model or not. A more complete analysis (Faul and Tipping, 2002), based on the second derivatives of the marginal likelihood, confirms these solutions are indeed the unique maxima of $\lambda(\alpha_i)$.

练习二阶导非负即归

Exercise 7.16

基于上节分析和推导，发展出的更高效的方法
超参数优化算法

Note that this approach has yielded a closed-form solution for α_i , for given values of the other hyperparameters. As well as providing insight into the origin of sparsity in the RVM, this analysis also leads to a practical algorithm for optimizing the hyperparameters that has significant speed advantages. This uses a fixed set of candidate basis vectors, and then cycles through them in turn to decide whether each vector should be included in the model or not. The resulting sequential sparse Bayesian learning algorithm is described below.

Sequential Sparse Bayesian Learning Algorithm

1. If solving a regression problem, initialize β .
2. Initialize using one basis function φ_1 , with hyperparameter α_1 set using (7.101), with the remaining hyperparameters α_j for $j \neq 1$ initialized to infinity, so that only φ_1 is included in the model.

3. Evaluate Σ and \mathbf{m} , along with q_i and s_i for all basis functions.
 4. Select a candidate basis function φ_i .
 5. If $q_i^2 > s_i$, and $\alpha_i < \infty$, so that the basis vector φ_i is already included in the model, then update α_i using (7.101).
 6. If $q_i^2 > s_i$, and $\alpha_i = \infty$, then add φ_i to the model, and evaluate hyperparameter α_i using (7.101).
 7. If $q_i^2 \leq s_i$, and $\alpha_i < \infty$ then remove basis function φ_i from the model, and set $\alpha_i = \infty$.
 8. If solving a regression problem, update β .
 9. If converged terminate, otherwise go to 3.

Note that if $q_i^2 \leq s_i$ and $\alpha_i = \infty$, then the basis function φ_i is already excluded from the model and no action is required.

前面介绍了算法流 // In practice, it is convenient to evaluate the quantities

程，下面介绍算法实现过程中在计算上的优化。

$$Q_i = \varphi_i^T C^{-1} t \quad (7.102)$$

$$S_i = \varphi_i^T \mathbf{C}^{-1} \varphi_i. \quad (7.103)$$

The quality and sparseness variables can then be expressed in the form

$$q_i = \frac{\alpha_i Q_i}{\alpha_i - S_i} \quad \text{这里的两个分子未被归一化} \quad (7.104)$$

$$s_i = \frac{\alpha_i S_i}{\alpha_i - S_i}. \quad (7.105)$$

Exercise 7.17

Note that when $\alpha_i = \infty$, we have $q_i = Q_i$ and $s_i = S_i$. Using (C.7), we can write

由此实际计算中先计算 $\left\{\begin{array}{l} (1.106) \\ [1.12] \end{array}\right.$

$$Q_i = \beta \varphi_i^T \mathbf{t} - \beta^2 \varphi_i^T \Phi \Sigma \Phi^T \mathbf{t} \quad (7.106)$$

$$S_i = \beta \varphi_i^T \varphi_i - \beta^2 \varphi_i^T \Phi \Sigma \Phi^T \varphi_i \quad (7.107)$$

where Φ and Σ involve only those basis vectors that correspond to finite hyperparameters α_i . At each stage the required computations therefore scale like $O(M^3)$, where M is the number of active basis vectors in the model and is typically much smaller than the number N of training patterns.

7.2.3 BVM for classification

We can extend the relevance vector machine framework to classification problems by applying the ARD prior over weights to a probabilistic linear classification model of the kind studied in Chapter 4. To start with, we consider two-class problems with a binary target variable $t \in \{0, 1\}$. The model now takes the form of a linear combination of basis functions transformed by a logistic sigmoid function

$$y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \phi(\mathbf{x})) \quad (7.108)$$

where $\sigma(\cdot)$ is the logistic sigmoid function defined by (4.59). If we introduce a Gaussian prior over the weight vector w , then we obtain the model that has been considered already in Chapter 4. The difference here is that in the RVM, this model uses the ARD prior (7.80) in which there is a separate precision hyperparameter associated with each weight parameter.

In contrast to the regression model, we can no longer integrate analytically over the parameter vector w . Here we follow Tipping (2001) and use the Laplace approximation, which was applied to the closely related problem of Bayesian logistic regression in Section 4.5.1.

Section 4.4

超参数学习的思路 (注意,
超参数的每项逻辑也会完成
对参数w的寻优,类似BSP(E))

对参数w的后验分布
 $p(w|t, \alpha)$ 和 Laplace
approximation

We begin by initializing the hyperparameter vector α . For this given value of α , we then build a Gaussian approximation to the posterior distribution and thereby obtain an approximation to the marginal likelihood. Maximization of this approximate marginal likelihood then leads to a re-estimated value for α , and the process is repeated until convergence.

[Let us consider the Laplace approximation for this model in more detail. For a fixed value of α , the mode of the posterior distribution over w is obtained by maximizing

$$\begin{aligned} \ln p(w|t, \alpha) &= \ln \{p(t|w)p(w|\alpha)\} - \ln p(t|\alpha) \\ &= \sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} - \frac{1}{2} w^T A w + \text{const} \end{aligned} \quad (7.109)$$

Exercise 7.18

where $A = \text{diag}(\alpha_i)$. This can be done using iterative reweighted least squares (IRLS) as discussed in Section 4.3.3. For this, we need the gradient vector and Hessian matrix of the log posterior distribution, which from (7.109) are given by

**注意,这里是针对
求导而不是对超参数**

$$\nabla \ln p(w|t, \alpha) = \Phi^T(t - y) - Aw \quad (7.110)$$

$$\nabla \nabla \ln p(w|t, \alpha) = -(\Phi^T B \Phi + A) \quad (7.111)$$

where B is an $N \times N$ diagonal matrix with elements $b_n = y_n(1-y_n)$, the vector $y = (y_1, \dots, y_N)^T$, and Φ is the design matrix with elements $\Phi_{ni} = \phi_i(x_n)$. Here we have used the property (4.88) for the derivative of the logistic sigmoid function. At convergence of the IRLS algorithm, the negative Hessian represents the inverse covariance matrix for the Gaussian approximation to the posterior distribution.

The mode of the resulting approximation to the posterior distribution, corresponding to the mean of the Gaussian approximation, is obtained setting (7.110) to zero, giving the mean and covariance of the Laplace approximation in the form

$$w^* = A^{-1} \Phi^T (t - y) \quad (7.112)$$

$$\Sigma = (\Phi^T B \Phi + A)^{-1}. \quad (7.113)$$

model evidence test [We can now use this Laplace approximation to evaluate the marginal likelihood. **通过超参数又可** Using the general result (4.135) for an integral evaluated using the Laplace approxi-

化 这里并不是直接利用 (4.135) 那么直接。我们要如何计算 $P(t|d) = \int p(t, w|d) dw$ 积分函数为 $p(t, w|d) \triangleq f(w)$ (积分变量只有 w , 不包含 t) , 我们应对 $\ln f(w)$ 进行二阶泰勒近似, 又 $p(t, w|d) = p(t|d) p(w|t, d)$, $p(t|d)$ 中不含 w

由式(7.13)得 $I_{\eta}(\mathbf{t}, \mathbf{w}) = I_{\eta} p(\mathbf{t}, \mathbf{w} | \alpha)$ 与 $I_{\eta} p(\mathbf{w} | \mathbf{t}, \alpha)$ 的情况一样，实际上就是对 $p(\mathbf{w} | \mathbf{t}, \alpha)$ 在进行 Laplace 近似，因此：

7.2. Relevance Vector Machines 355

$$p(\mathbf{t} | \alpha) = \int p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \mathbf{t}, \alpha) d\mathbf{w} = p(\mathbf{t} | \alpha) \int p(\mathbf{w} | \mathbf{t}, \alpha) d\mathbf{w}$$

mation, we have 由式(4.13丁)得 $\approx p(\mathbf{t} | \alpha) p(\mathbf{w}^* | \mathbf{t}, \alpha) (2\pi)^{M/2} |\Sigma|^{1/2}$

$$\begin{aligned} p(\mathbf{t} | \alpha) &= \int p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \alpha) d\mathbf{w} = p(\mathbf{t} | \alpha) p(\mathbf{w}^* | \mathbf{t}, \alpha) (2\pi)^{M/2} |\Sigma|^{1/2} \\ &= p(\mathbf{t} | \mathbf{w}^*) p(\mathbf{w}^* | \alpha) (2\pi)^{M/2} |\Sigma|^{1/2} \\ &\simeq p(\mathbf{t} | \mathbf{w}^*) p(\mathbf{w}^* | \alpha) (2\pi)^{M/2} |\Sigma|^{1/2} \end{aligned} \quad (7.114)$$

注意该推导和 RVM 中 KIC 的推导类似。

If we substitute for $p(\mathbf{t} | \mathbf{w}^*)$ and $p(\mathbf{w}^* | \alpha)$ and then set the derivative of the marginal likelihood with respect to α_i equal to zero, we obtain

$$-\frac{1}{2}(w_i^*)^2 + \frac{1}{2\alpha_i} - \frac{1}{2}\sum_{ii} = 0. \quad (7.115)$$

Defining $\gamma_i = 1 - \alpha_i \sum_{ii}$ and rearranging then gives

$$\alpha_i^{\text{new}} = \frac{\gamma_i}{(w_i^*)^2} \quad (7.116)$$

which is identical to the re-estimation formula (7.87) obtained for the regression RVM.

// If we define

$$\hat{\mathbf{t}} = \Phi \mathbf{w}^* + \mathbf{B}^{-1}(\mathbf{t} - \mathbf{y}) \quad (7.117)$$

we can write the approximate log marginal likelihood in the form

$$\ln p(\mathbf{t} | \alpha) = -\frac{1}{2} \left\{ N \ln(2\pi) + \ln |\mathbf{C}| + (\hat{\mathbf{t}})^T \mathbf{C}^{-1} \hat{\mathbf{t}} \right\} \quad (7.118)$$

where

$$\mathbf{C} = \mathbf{B} + \Phi \mathbf{A} \Phi^T. \quad (7.119)$$

This takes the same form as (7.85) in the regression case, and so we can apply the same analysis of sparsity and obtain the same fast learning algorithm in which we fully optimize a single hyperparameter α_i at each step.]

SSVM in RVM Figure 7.12 shows the relevance vector machine applied to a synthetic classification data set. [

We see that the relevance vectors tend not to lie in the region of the decision boundary, in contrast to the support vector machine. This is consistent with our earlier discussion of sparsity in the RVM, because a basis function $\phi_i(\mathbf{x})$ centred on a data point near the boundary will have a vector φ_i that is poorly aligned with the training data vector \mathbf{t} . [这句话的逻辑是：RVM 中一般是 localized，通常越靠近中心点处 $y = \mathbf{w}^T \phi(\mathbf{x})$ ，则不存在 relevance vectors 的概念。 One of the potential advantages of the relevance vector machine compared with the SVM is that it makes probabilistic predictions. For example, this allows the RVM to be used to help construct an emission density in a nonlinear extension of the linear dynamical system for tracking faces in video sequences (Williams *et al.*, 2005).]

Appendix A

注意，这里的例子

使用的是 SVM-like 模型

即式(7.78)，若使用一般地

$y = \mathbf{w}^T \phi(\mathbf{x})$ ，则不存在 relevance

vectors 的概念。

Section 13.3

分类的情况 So far, we have considered the RVM for binary classification problems. For $K > 2$ classes, we again make use of the probabilistic approach in Section 4.3.4 in which there are K linear models of the form

$$\phi_i(\mathbf{x}) \text{ 的值越大, 越靠近边界, 则说明越靠近 } a_k = \mathbf{w}_k^T \mathbf{x} \quad (7.120)$$

边界中 $\phi_i(\mathbf{x})$ 的值越大, 这与间隔 $y_i(\mathbf{x})$ 的离散情况并不

匹配, 边界附近训练点会使 $y_i(\mathbf{x})$ 接近 0, 又 $y_i(\mathbf{x})$ 实现了

对训练数据较好地分离, 则 $y_i(\mathbf{x})$ 对训练数据的分离会较差, 即 ϕ_i is poorly aligned with t .

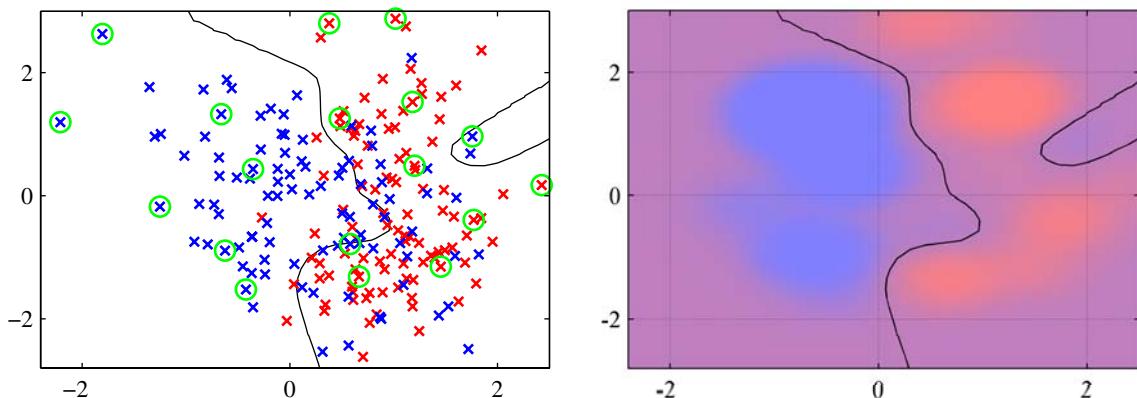


Figure 7.12 Example of the relevance vector machine applied to a synthetic data set, in which the left-hand plot shows the decision boundary and the data points, with the relevance vectors indicated by circles. Comparison with the results shown in Figure 7.4 for the corresponding support vector machine shows that the RVM gives a much sparser model. The right-hand plot shows the posterior probability given by the RVM output in which the proportion of red (blue) ink indicates the probability of that point belonging to the red (blue) class.

which are combined using a softmax function to give outputs

$$y_k(\mathbf{x}) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}. \quad (7.121)$$

The log likelihood function is then given by

$$\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad (7.122)$$

where the target values t_{nk} have a 1-of- K coding for each data point n , and \mathbf{T} is a matrix with elements t_{nk} . Again, the Laplace approximation can be used to optimize the hyperparameters (Tipping, 2001), in which the model and its Hessian are found using IRLS. This gives a more principled approach to multiclass classification than the pairwise method used in the support vector machine and also provides probabilistic predictions for new data points. The principal disadvantage is that the Hessian matrix has size $MK \times MK$, where M is the number of active basis functions, which gives an additional factor of K^3 in the computational cost of training compared with the two-class RVM.

The principal disadvantage of the relevance vector machine is the relatively long training times compared with the SVM. This is offset, however, by the avoidance of cross-validation runs to set the model complexity parameters. Furthermore, because it yields sparser models, the computation time on test points, which is usually the more important consideration in practice, is typically much less.]

前面介紹過
時時機模型及訓練方法
后面討論了SVM的優劣，
屬於是全面碾壓SVM了

Exercises

- 7.1** (**) **www** Suppose we have a data set of input vectors $\{\mathbf{x}_n\}$ with corresponding target values $t_n \in \{-1, 1\}$, and suppose that we model the density of input vectors within each class separately using a Parzen kernel density estimator (see Section 2.5.1) with a kernel $k(\mathbf{x}, \mathbf{x}')$. Write down the minimum misclassification-rate decision rule assuming the two classes have equal prior probability. Show also that, if the kernel is chosen to be $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$, then the classification rule reduces to simply assigning a new input vector to the class having the closest mean. Finally, show that, if the kernel takes the form $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$, that the classification is based on the closest mean in the feature space $\phi(\mathbf{x})$.
- 7.2** (*) Show that, if the 1 on the right-hand side of the constraint (7.5) is replaced by some arbitrary constant $\gamma > 0$, the solution for the maximum margin hyperplane is unchanged.
- 7.3** (**) Show that, irrespective of the dimensionality of the data space, a data set consisting of just two data points, one from each class, is sufficient to determine the location of the maximum-margin hyperplane.
- 7.4** (**) **www** Show that the value ρ of the margin for the maximum-margin hyperplane is given by

$$\frac{1}{\rho^2} = \sum_{n=1}^N a_n \quad (7.123)$$

where $\{a_n\}$ are given by maximizing (7.10) subject to the constraints (7.11) and (7.12).

- 7.5** (**) Show that the values of ρ and $\{a_n\}$ in the previous exercise also satisfy

$$\frac{1}{\rho^2} = 2\tilde{L}(\mathbf{a}) \quad (7.124)$$

where $\tilde{L}(\mathbf{a})$ is defined by (7.10). Similarly, show that

$$\frac{1}{\rho^2} = \|\mathbf{w}\|^2. \quad (7.125)$$

- 7.6** (*) Consider the logistic regression model with a target variable $t \in \{-1, 1\}$. If we define $p(t = 1|y) = \sigma(y)$ where $y(\mathbf{x})$ is given by (7.1), show that the negative log likelihood, with the addition of a quadratic regularization term, takes the form (7.47).
- 7.7** (*) Consider the Lagrangian (7.56) for the regression support vector machine. By setting the derivatives of the Lagrangian with respect to \mathbf{w} , b , ξ_n , and $\hat{\xi}_n$ to zero and then back substituting to eliminate the corresponding variables, show that the dual Lagrangian is given by (7.61).

- 7.8** (★) **www** For the regression support vector machine considered in Section 7.1.4, show that all training data points for which $\xi_n > 0$ will have $a_n = C$, and similarly all points for which $\hat{\xi}_n > 0$ will have $\hat{a}_n = C$.
- 7.9** (★) Verify the results (7.82) and (7.83) for the mean and covariance of the posterior distribution over weights in the regression RVM.
- 7.10** (★★) **www** Derive the result (7.85) for the marginal likelihood function in the regression RVM, by performing the Gaussian integral over \mathbf{w} in (7.84) using the technique of completing the square in the exponential.
- 7.11** (★★) Repeat the above exercise, but this time make use of the general result (2.115).
- 7.12** (★★) **www** Show that direct maximization of the log marginal likelihood (7.85) for the regression relevance vector machine leads to the re-estimation equations (7.87) and (7.88) where γ_i is defined by (7.89).
- 7.13** (★★) In the evidence framework for RVM regression, we obtained the re-estimation formulae (7.87) and (7.88) by maximizing the marginal likelihood given by (7.85). Extend this approach by inclusion of hyperpriors given by gamma distributions of the form (B.26) and obtain the corresponding re-estimation formulae for α and β by maximizing the corresponding posterior probability $p(\mathbf{t}, \alpha, \beta | \mathbf{X})$ with respect to α and β .
- 7.14** (★★) Derive the result (7.90) for the predictive distribution in the relevance vector machine for regression. Show that the predictive variance is given by (7.91).
- 7.15** (★★) **www** Using the results (7.94) and (7.95), show that the marginal likelihood (7.85) can be written in the form (7.96), where $\lambda(\alpha_n)$ is defined by (7.97) and the sparsity and quality factors are defined by (7.98) and (7.99), respectively.
- 7.16** (★) By taking the second derivative of the log marginal likelihood (7.97) for the regression RVM with respect to the hyperparameter α_i , show that the stationary point given by (7.101) is a maximum of the marginal likelihood.
- 7.17** (★★) Using (7.83) and (7.86), together with the matrix identity (C.7), show that the quantities S_n and Q_n defined by (7.102) and (7.103) can be written in the form (7.106) and (7.107).
- 7.18** (★) **www** Show that the gradient vector and Hessian matrix of the log posterior distribution (7.109) for the classification relevance vector machine are given by (7.110) and (7.111).
- 7.19** (★★) Verify that maximization of the approximate log marginal likelihood function (7.114) for the classification relevance vector machine leads to the result (7.116) for re-estimation of the hyperparameters.