



MACHINE LEARNING

CHAPTER 4: KERNEL METHODS

KERNEL METHODS

1. Kernel function and dual representation
2. Construction kernel
3. Gaussian process regression

Kernel function

For a nonlinear feature space mapping $\phi(\mathbf{x})$, the *kernel function* is given by

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

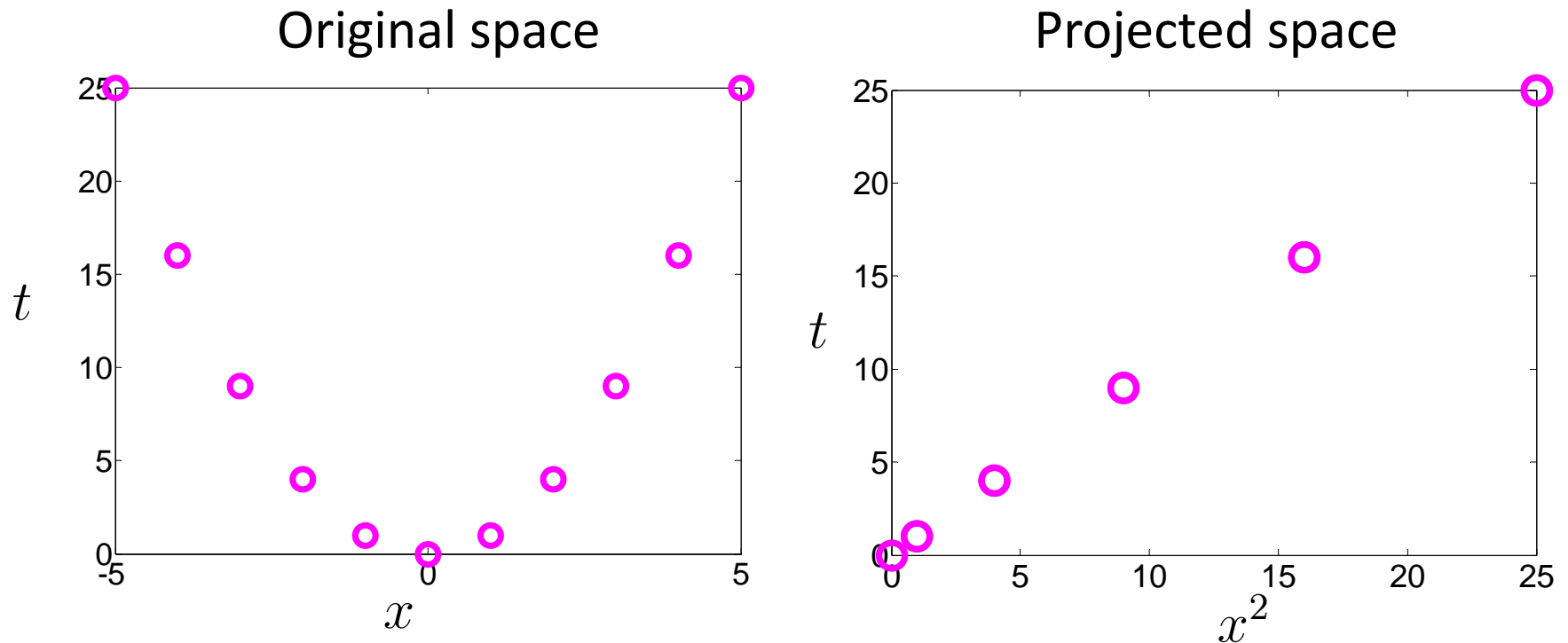
From this definition, we see that the kernel is a symmetric function.

Dual representation

A linear regression model whose parameters are determined by minimizing a regularized sum-of-squares error function given by

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

Why do we use kernel trick



Basis function: $\phi(x) = x^2$,
linear regression model: $y = \phi(x)$.

Why do we use kernel trick

If the projected feature has infinite dimensionality, it is impossible to learn the regression parameters \mathbf{w} .

But the kernel function under this situation has a simple expression, such as Gaussian kernel

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

Dual problem is much more easy.

Dual representation

We suppose

$$\mathbf{w} = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \mathbf{\Phi}^T \mathbf{a}$$

where

$$\mathbf{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}.$$

Dual representation

Substitute $\mathbf{w} = \Phi^T \mathbf{a}$ into $J(\mathbf{w})$, we obtain

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

We define the Gram matrix $\mathbf{K} = \Phi \Phi^T$, which is an $N \times N$ symmetric matrix with elements

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

Dual representation

In terms of Gram matrix, the sum-of-squares error function can be rewritten as

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

Setting the gradient of $J(\mathbf{a})$ with respect to zero, we obtain the following solution

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

Dual representation

If we substitute this solution back to the linear regression model, the prediction function is

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

Where $\mathbf{k}(\mathbf{x})$ is an N dimensional vector with elements

$$k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$$

KERNEL METHODS

1. Kernel function and dual representation
2. Construction kernels
3. Gaussian process regression

Constructing kernels

a function $k(\mathbf{x}, \mathbf{x}')$ is a valid kernel
if 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 $\{\mathbf{x}_n\}$

Polynomial kernel $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^2$

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

Techniques for constructing new kernels

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

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

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

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

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

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

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

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

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

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

KERNEL METHODS

1. Kernel function and dual representation
2. Construction kernels
3. Gaussian process regression

Gaussian process regression

A general regression model

$$t_n = y_n + \epsilon_n$$

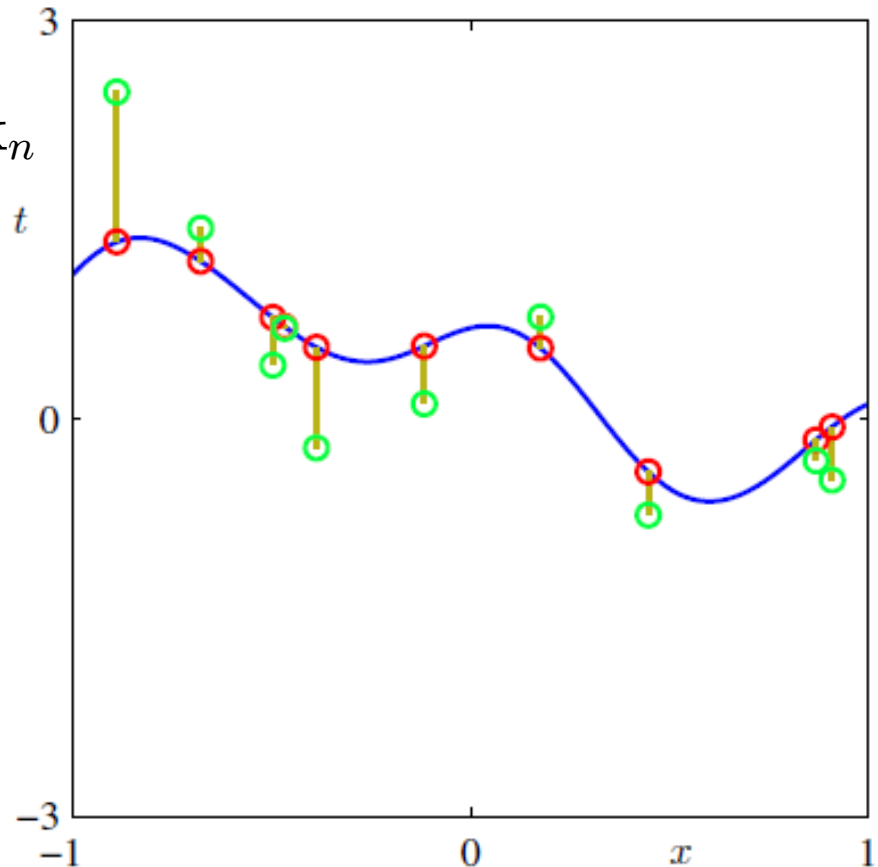
where $y_n = y(\mathbf{x}_n)$ is the output of regression model, and ϵ_n is a random noise variable. We suppose ϵ_n has a Gaussian distribution, so that

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

Gaussian process regression

y_n is the output of the prediction function $y(\mathbf{x})$ at \mathbf{x}_n

$$t_n = y_n + \epsilon_n$$



Gaussian process regression

The joint distribution of target values

$\mathbf{t} = (t_1, \dots, t_N)^T$ conditioned on the values of \mathbf{y} is

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

Gaussian process regression

Gaussian process regression suppose a Gaussian process prior on the outputs of regression model $\mathbf{y} = [y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)]^T$

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

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 strongly correlated than for dissimilar points.

Gaussian process regression

$$\text{cov} (y (\mathbf{x}_n) , y (\mathbf{x}_m)) = k (\mathbf{x}_n, \mathbf{x}_m)$$

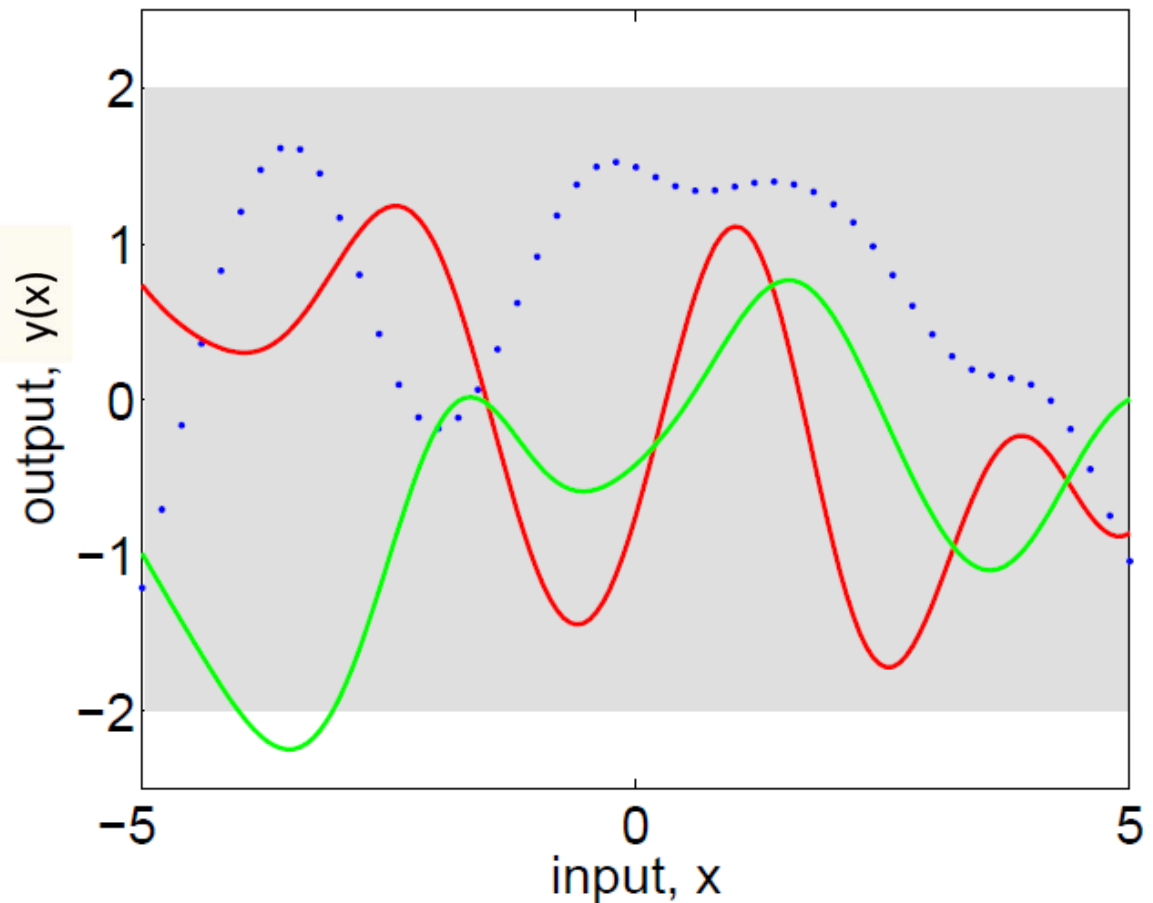
One widely used kernel function for Gaussian process is given by the exponential of a quadratic form,

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

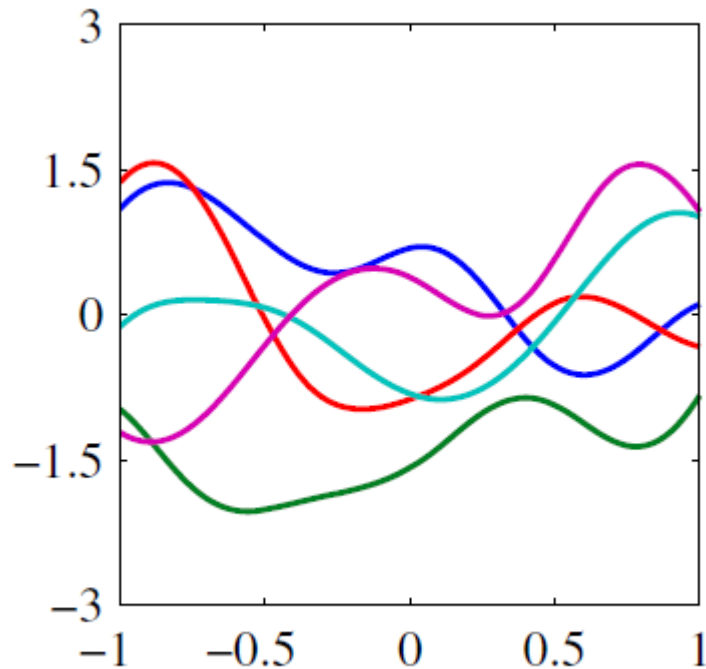
here $\theta = \{\theta_0, \theta_1, \theta_2, \theta_3, \beta\}$

Samples from Gaussian process prior

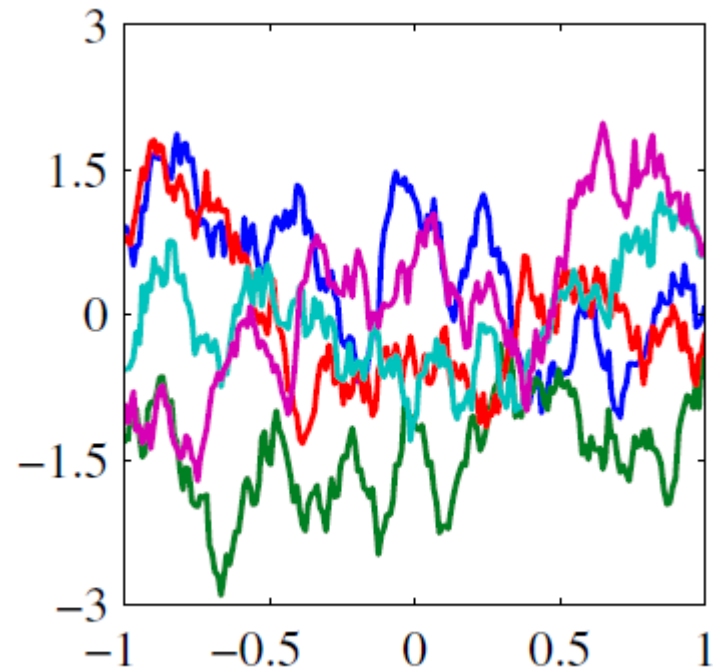
Three functions drawn at random from a Gaussian process prior, the dots indicate values of y actually generated; the two other functions have been drawn as lines by joining a large number of evaluated points



Samples from Gaussian process prior



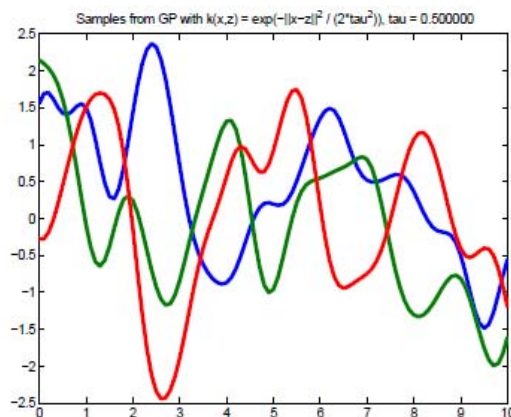
Gaussian kernel



Exponential kernel

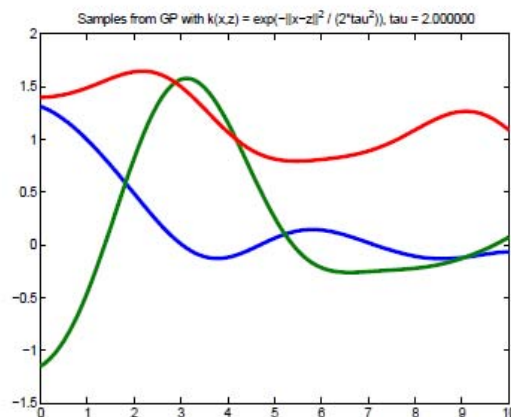
Samples from Gaussian process prior

Samples from a zero mean Gaussian process prior with covariance function of Gaussian kernel



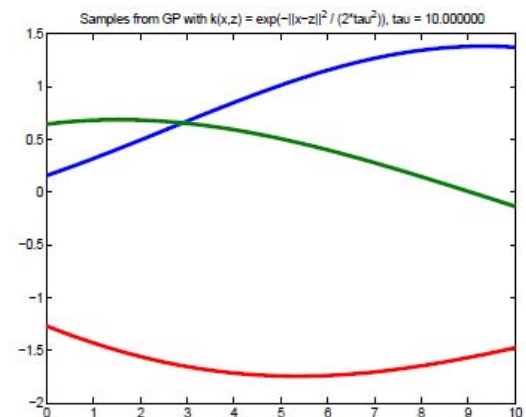
(a)

$$\sigma = 0.5$$



(b)

$$\sigma = 2$$



(c)

$$\sigma = 10$$

Samples from Gaussian process prior

As we know the Gaussian kernel is

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\|\mathbf{x}_n - \mathbf{x}_m\|^2 / 2\sigma^2)$$

The larger the hyper-parameter σ is, the more strongly correlated $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ are.

Gaussian process regression

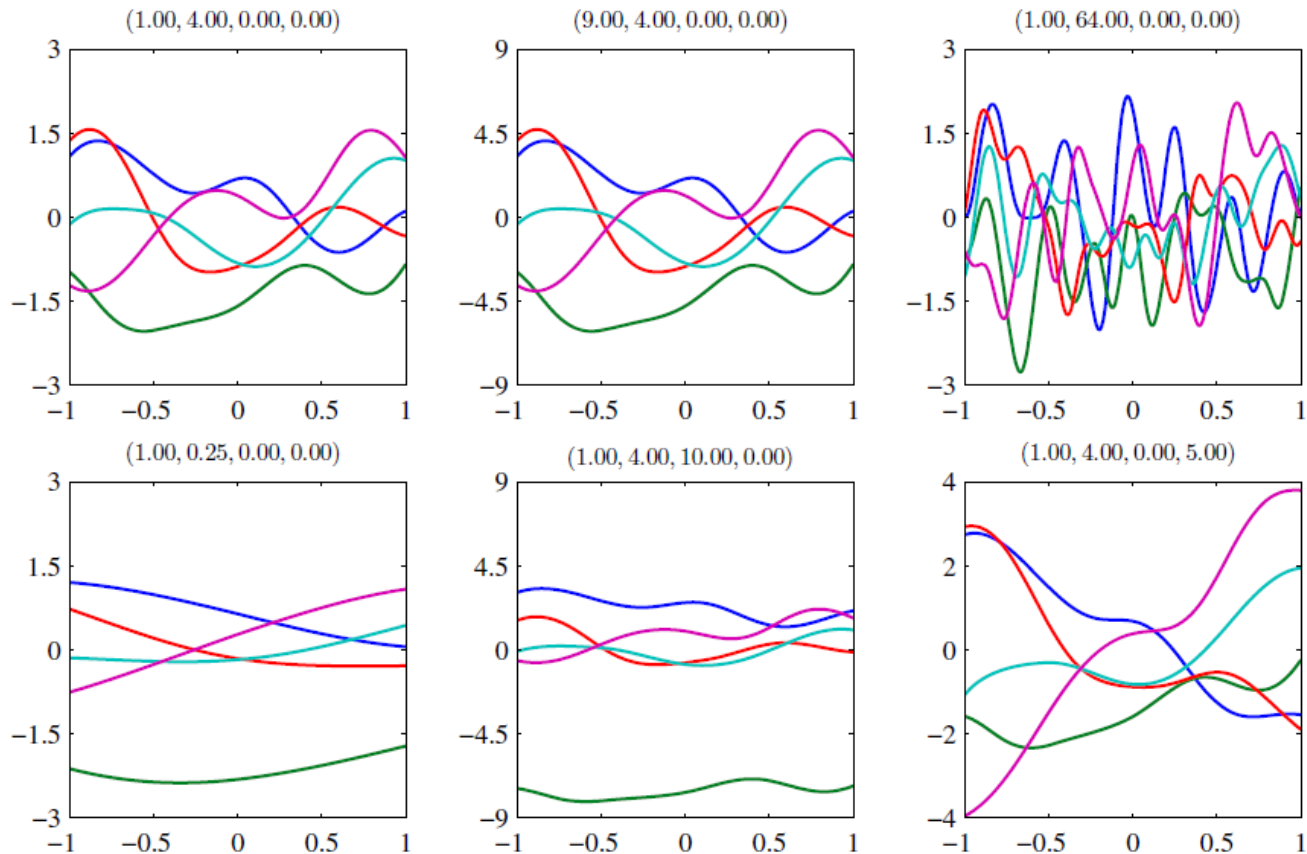


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

Gaussian process regression

The marginal distribution of target values

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

where the covariance matrix \mathbf{C} has elements

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

Gaussian process regression

The objective function of Gaussian process regression

$$\ln p(\mathbf{t}|\mathbf{X}) = -\frac{1}{2} \ln |\mathbf{C}| - \frac{1}{2} \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} - \frac{N}{2} \ln 2\pi$$

It is nonconvex.

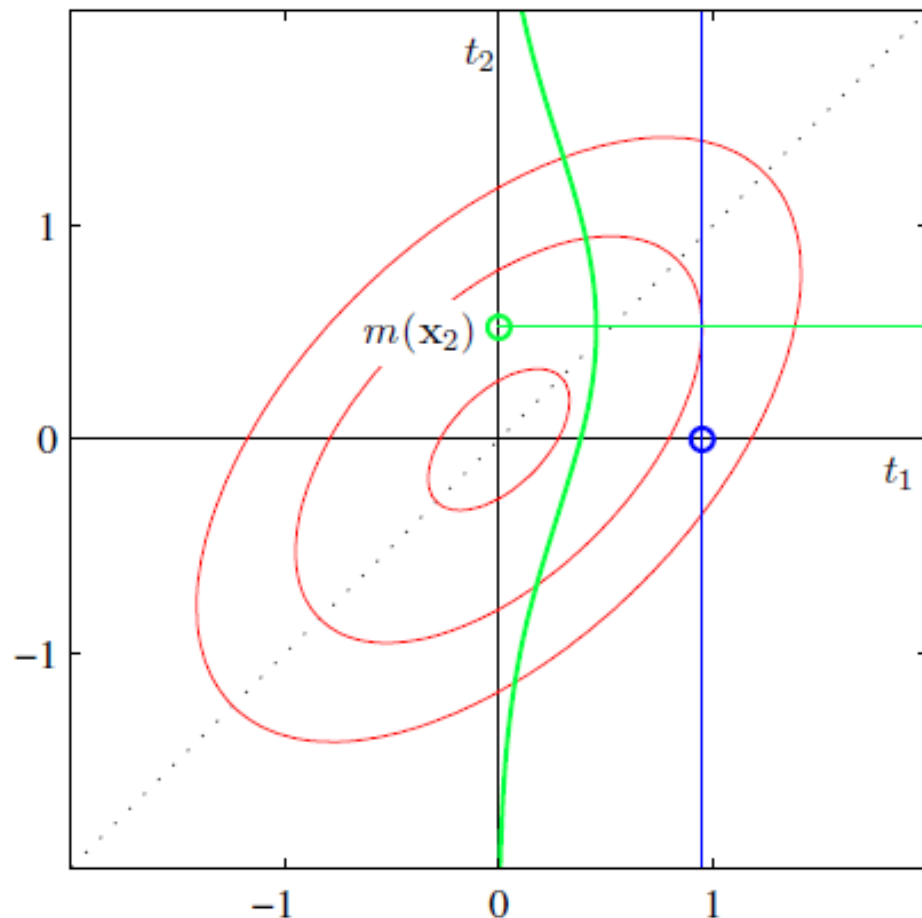
Gaussian process regression

To get the optimal hyper-parameters, we use gradient descent method.

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

Gaussian process regression

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.



Gaussian process regression

Given a new input vector \mathbf{x}_{N+1} , the joint distribution of target values is

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

where the covariance matrix is

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

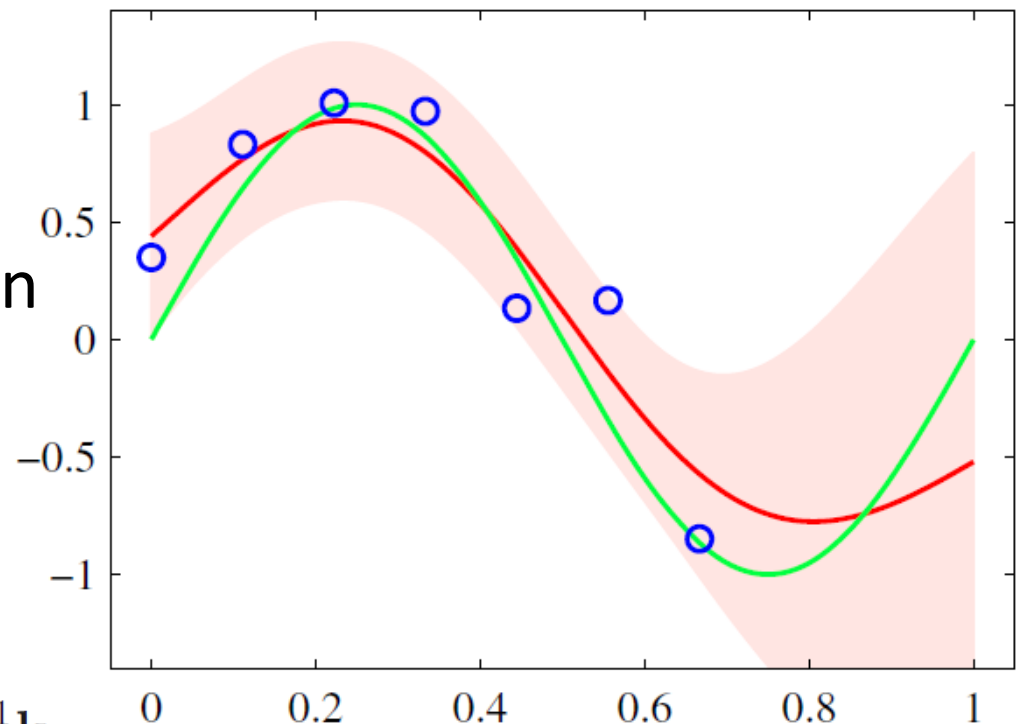
Gaussian process regression

The conditional distribution

$$p(t_{N+1} | \mathbf{X}, \mathbf{t}, \mathbf{x}_{N+1})$$

is a Gaussian
Distribution with mean
and covariance given
by

$$\begin{aligned} m(\mathbf{x}_{N+1}) &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \\ \sigma^2(\mathbf{x}_{N+1}) &= c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \end{aligned}$$



Gaussian process regression

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

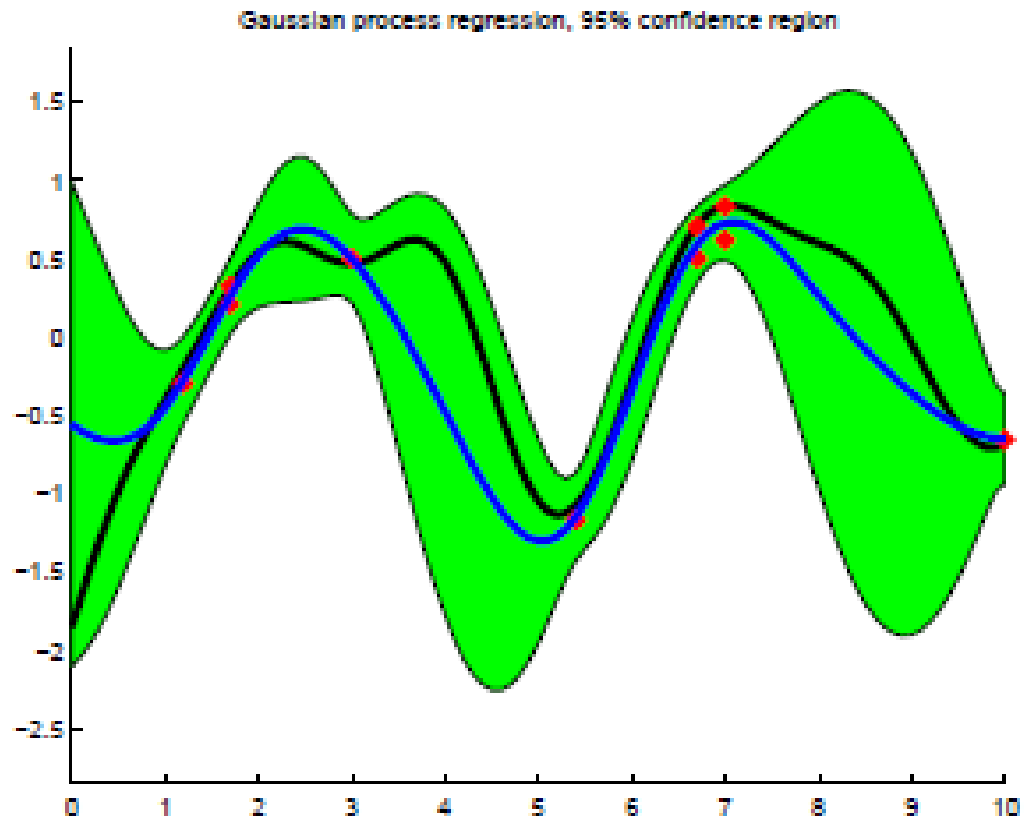
The mean of the predictive distribution can be written as a function of \mathbf{x}_{N+1}

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

where a_n is the n th component of $\mathbf{C}_N^{-1} \mathbf{t}$.

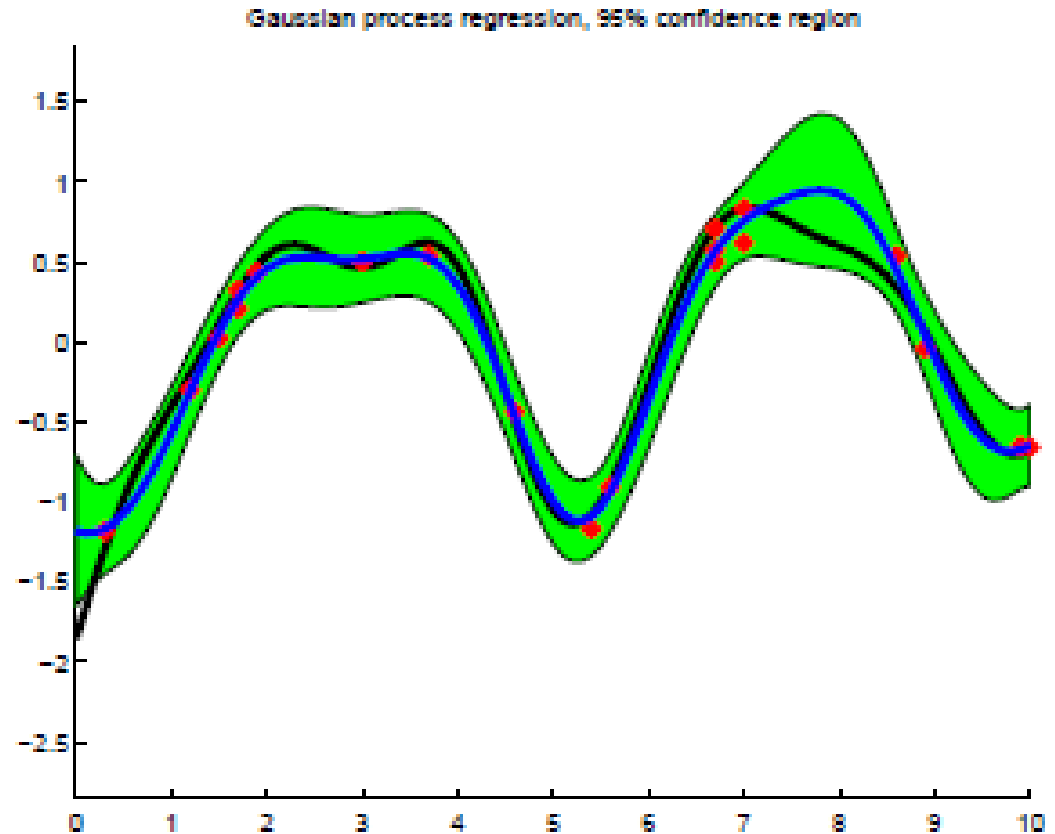
Gaussian process regression

The number of training examples is 10, noise level is $\sigma = 1$



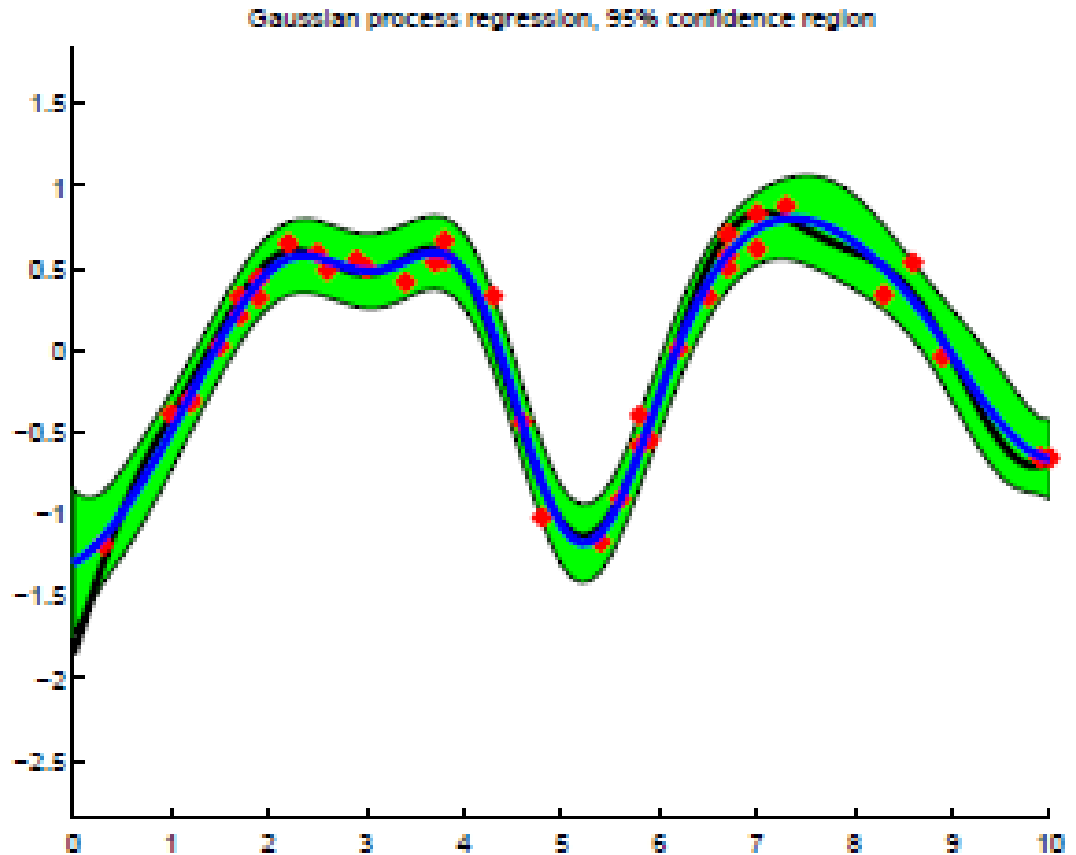
Gaussian process regression

The number of
training examples
is 20, noise level
is $\sigma = 1$



Gaussian process regression

The number of
training examples
is 40, noise level
is $\sigma = 1$



Linear regression revisited

Consider the linear model

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

A prior distribution over \mathbf{w} given by an isotropic Gaussian of the form

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

governed by the hyper-parameter α .

Linear regression revisited

Let \mathbf{y} denote model outputs $\mathbf{y} = [y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)]^T$,
we have

$$\mathbf{y} = \Phi \mathbf{w}$$

where

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

Linear regression revisited

\mathbf{y} is a linear transformation of \mathbf{x} , therefore \mathbf{y} is Gaussian.

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

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

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

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

Linear regression revisited

This model provides us with a particular example of Gaussian process.

The set of values of $y(\mathbf{x})$ evaluated at an arbitrary set of point $\mathbf{x}_1, \dots, \mathbf{x}_N$ jointly have a Gaussian distribution.
