



Pattern Recognition and Machine Learning

Chapter 6: Kernel Methods

Polynomial regression

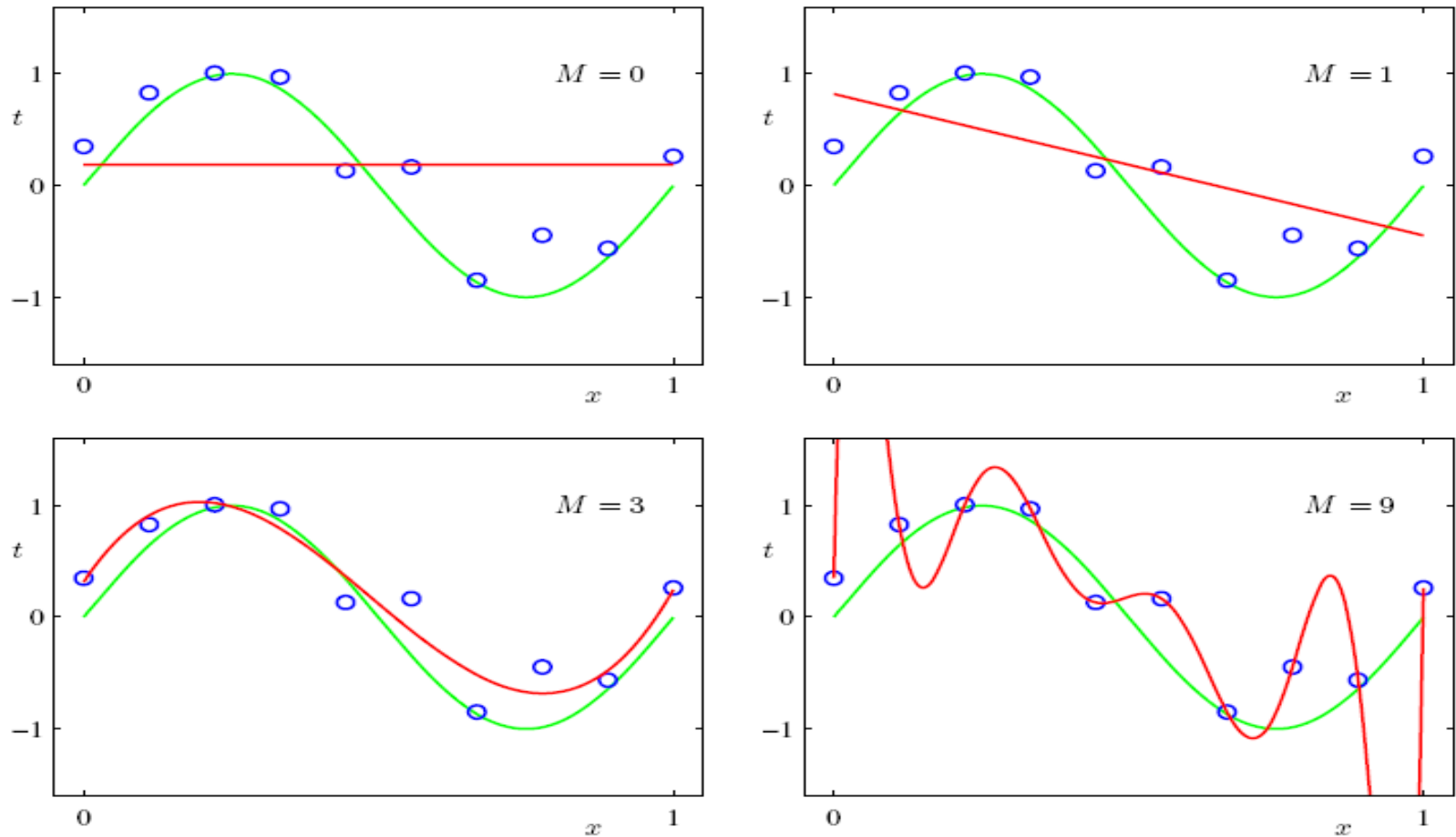
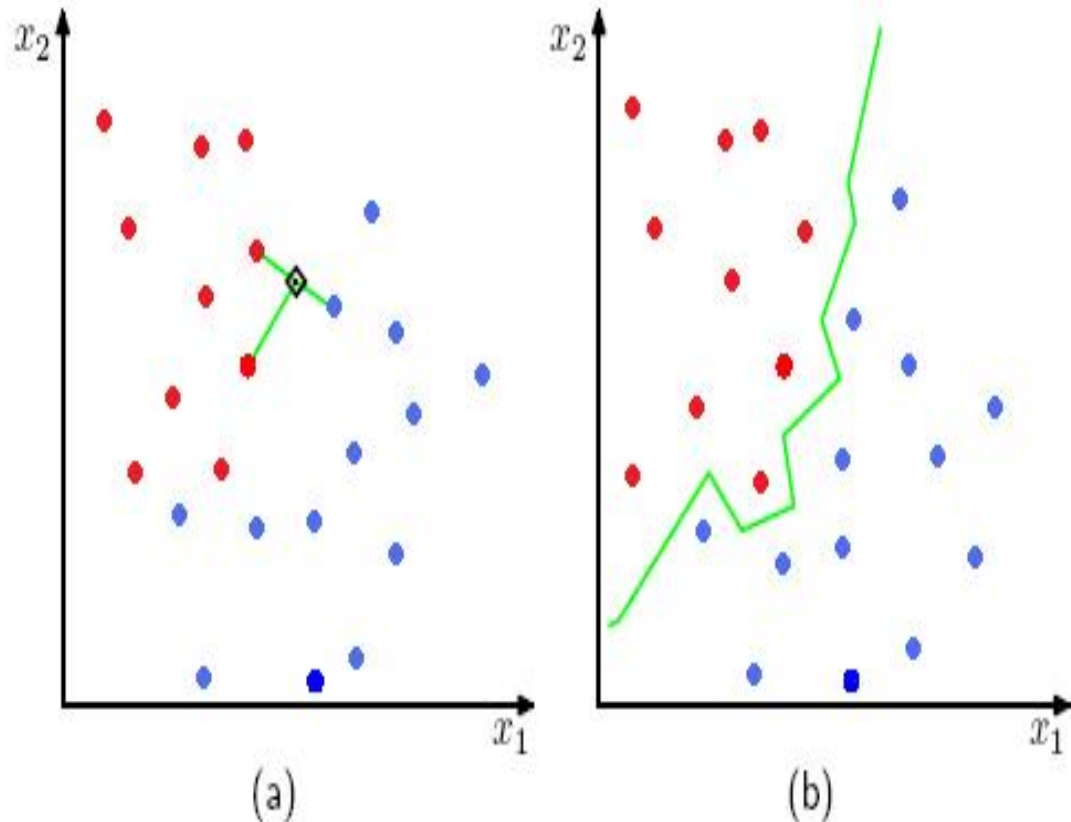


Figure 1.4 Plots of polynomials having various orders M , shown as red curves, fitted to the data set shown in Figure 1.2.

K-Nearest Neighbour

Figure 2.27 (a) In the K -nearest-neighbour classifier, a new point, shown by the black diamond, is classified according to the majority class membership of the K closest training data points, in this case $K = 3$. (b) In the nearest-neighbour ($K = 1$) approach to classification, the resulting decision boundary is composed of hyperplanes that form perpendicular bisectors of pairs of points from different classes.



Kernel Methods

- Are used to find and study general types of relations (eg. Clustering, correlation, classification) in general types of data (vectors, images, sets).
- For nonlinear feature space mapping $\phi(\mathbf{x})$, the kernel function is given as

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

Kernel Methods

- Kernel Trick (Kernel Substitution) :- If the input vector enters in the form of scalar product, then the scalar product can be replaced by some other kernel functions .

- Linear Kernels: $k(\mathbf{X}, \mathbf{X}') = \mathbf{x}^T \mathbf{x}'$

- Stationary Kernels: $k(\mathbf{X}, \mathbf{X}') = k(\mathbf{X} - \mathbf{X}')$

- Homogeneous Kernels : $k(\mathbf{X}, \mathbf{X}') = k(\|\mathbf{X} - \mathbf{X}'\|)$.

Dual Representations

- Consider linear regression model whose regularized sum-of-squares error function is given by

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad \text{where } \lambda \geq 0.$$

- Minimizing with respect to \mathbf{w} ,

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

where Φ is the design matrix, whose n th row is given by $\phi(\mathbf{x}_n)^T$.

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

$$\mathbf{a} = (a_1, \dots, a_N)^T$$

Dual Representations

- Substituting $\mathbf{w} = \Phi^T \mathbf{a}$ into $J(\mathbf{w})$

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

where $\mathbf{t} = (t_1, \dots, t_N)^T$.

- Define *Gram Matrix* $\mathbf{K} = \Phi \Phi^T$, which is an $N \times N$ symmetric matrix with elements

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

- Re-writing the sum-of-squares error using the *Gram Matrix*

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

Dual Representations

- Minimizing $J(\mathbf{a})$ with respect to \mathbf{a} , and solving

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

- Substituting back to the original linear regression model

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

- Note that vector \mathbf{a} is computed by inverting $N \times N$ matrix.
-

- *Compare:* $\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$

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

Dual Representations

- Even though the dual formulation requires higher dimensions , it comes with some advantages.
 - We can see that $y(\mathbf{x})$ is expressed entirely in terms of the kernel functions.
 - This allows us to work with feature spaces of higher dimension, even with infinite dimensionality.

Constructing Kernels

- First Approach: Select the feature space mapping $\phi(x)$ and use it to construct the corresponding kernel.

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

- Examples of basis function:

- Polynomial: $\phi_j(x) = x^j$

- Gaussian: $\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$

- Sigmoid: $\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right)$

Constructing Kernels

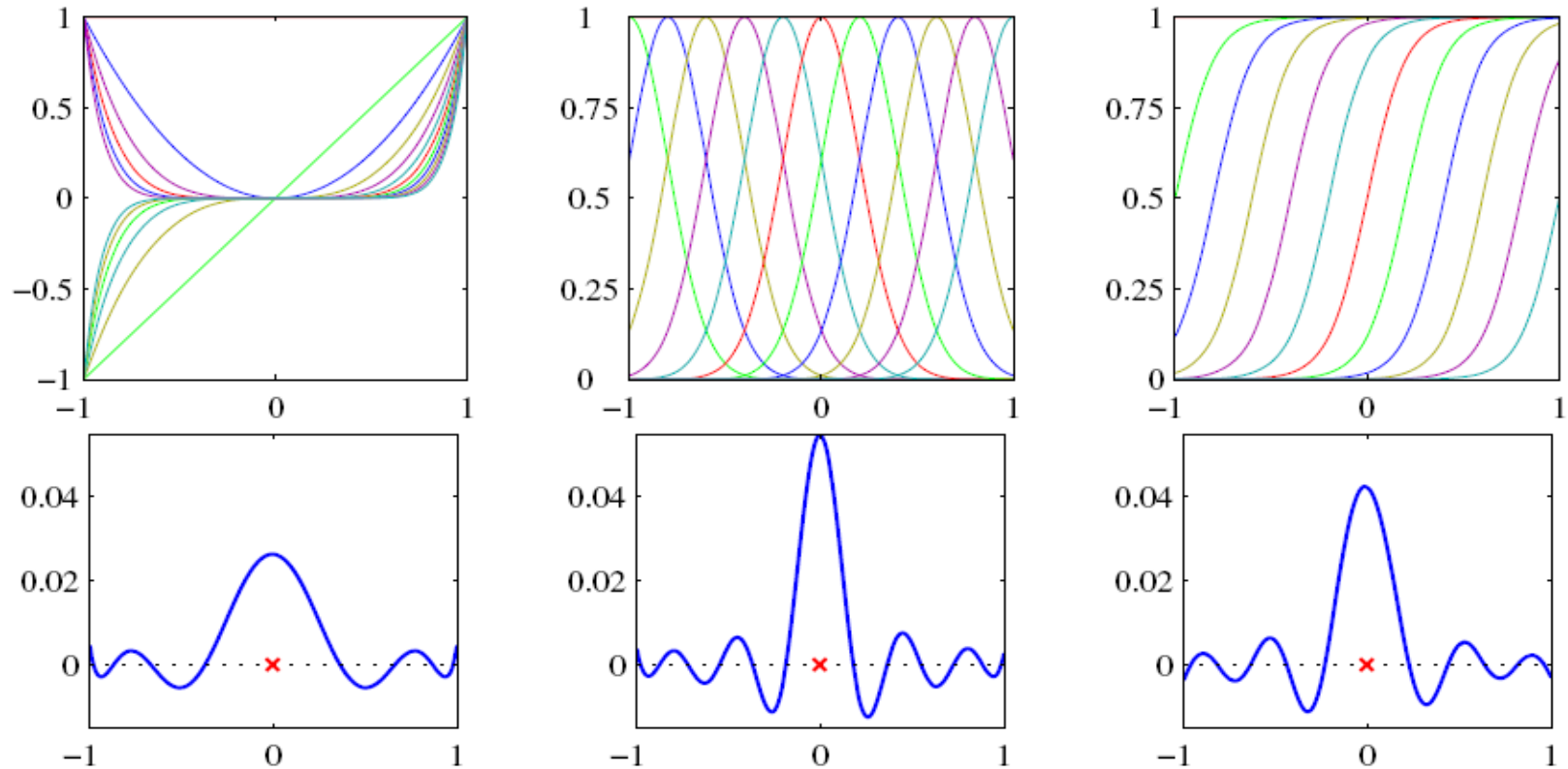


Figure 6.1 Illustration of the construction of kernel functions starting from a corresponding set of basis functions. In each column the lower plot shows the kernel function $k(x, x')$ defined by (6.10) plotted as a function of x for $x' = 0$, while the upper plot shows the corresponding basis functions given by polynomials (left column), 'Gaussians' (centre column), and logistic sigmoids (right column).

Constructing Kernels

- Second Approach: Construct kernel functions directly. We need to make sure that we are selecting a valid kernel.
- Valid Kernels: kernels who has a Gram Matrix whose components are **positive semi-definite** for all possible choices of input data.

Constructing Kernels

- Consider a simple example: $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$.
- Considering a 2D input space

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

- New complex Kernels can also be reconstructed by using simpler kernels as building blocks.

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.

Constructing Kernels

- Gaussian Kernels

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

Taking the inner part:

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

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

Substituting kernel

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

Linear Regression Revisited

- Consider

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

- Taking prior distribution over \mathbf{w}

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

- Note that this induces a probability distribution over function $y(\mathbf{x})$.

- In terms of vector representation

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

Linear Regression Revisited

- Note that since \mathbf{y} is a linear combination of Gaussian distributed values \mathbf{w} , and hence it is also a Gaussian distributed.

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

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

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

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

- Definition: A Gaussian process is defined as a probability distribution over functions $y(\mathbf{x})$ evaluated at an arbitrary set of points $\mathbf{x}_1, \dots, \mathbf{x}_N$ jointly have a Gaussian distribution.

Linear Regression Revisited

- Note that the joint distribution of Gaussian process can be completely specified by the mean and covariance.
- Note also that the covariance can be evaluated from the kernel function.

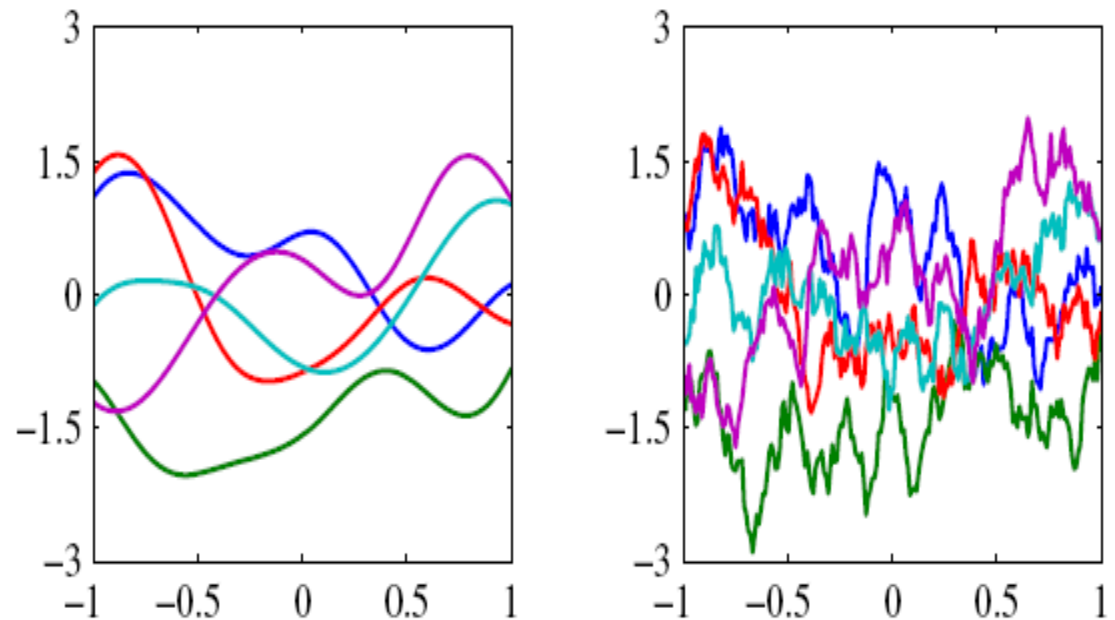
$$\mathbb{E} [y(\mathbf{x}_n)y(\mathbf{x}_m)] = k(\mathbf{x}_n, \mathbf{x}_m).$$

- Taking the Gaussian and exponential kernel functions, for example, samples are drawn.

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

Linear Regression Revisited

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



Gaussian processes for regression

- Considering noise on the observed target

$$t_n = y_n + \epsilon_n$$

- Considering Gaussian noise

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

- For independent noise, the joint distribution is given by

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

- Using the Gaussian process

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

- Marginalizing the probability

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

Gaussian processes for regression

- Where \mathbf{C} is the covariance matrix given as

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

- Note the summation of the covariance.
- Widely used Kernel function

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

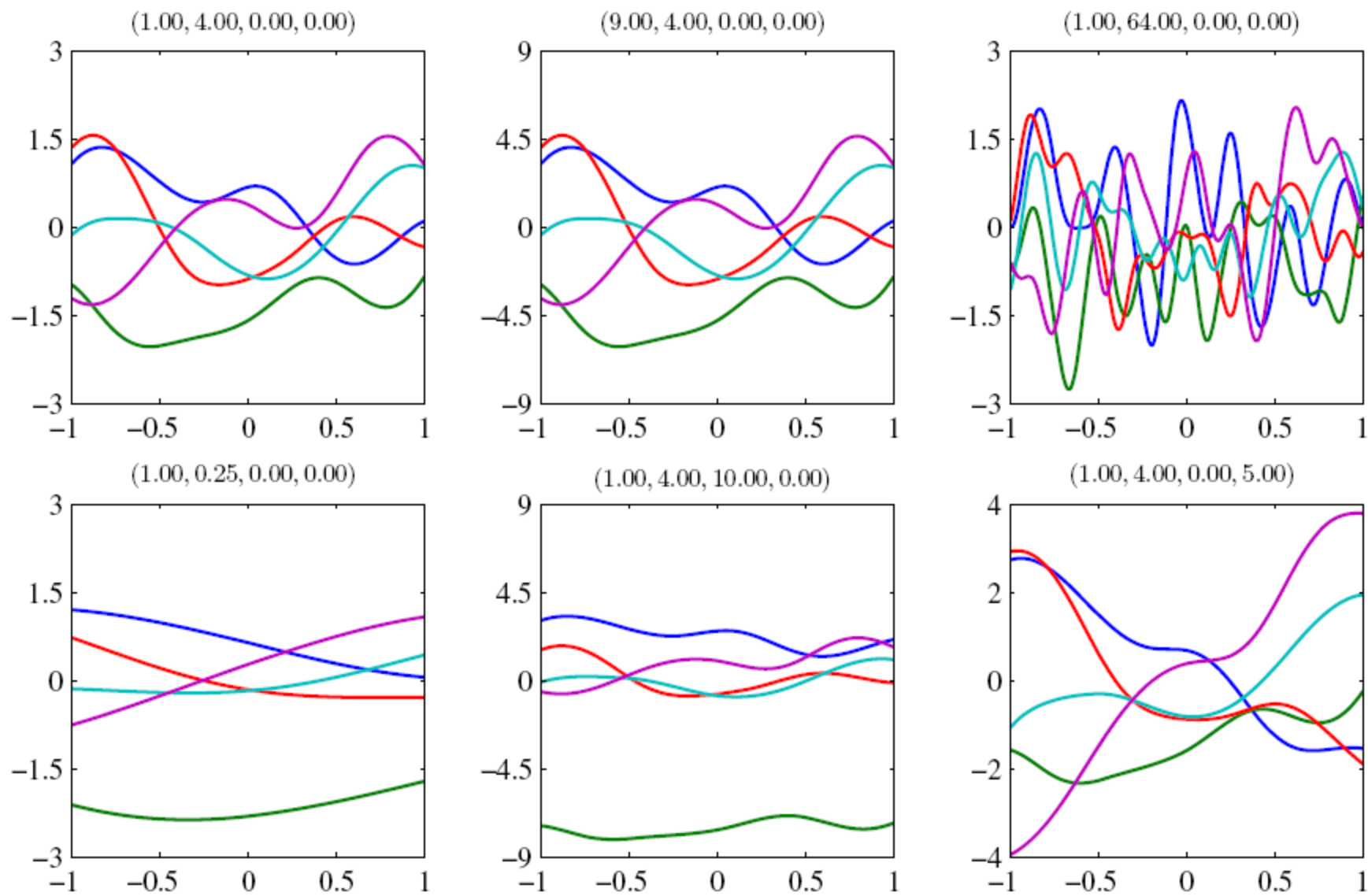
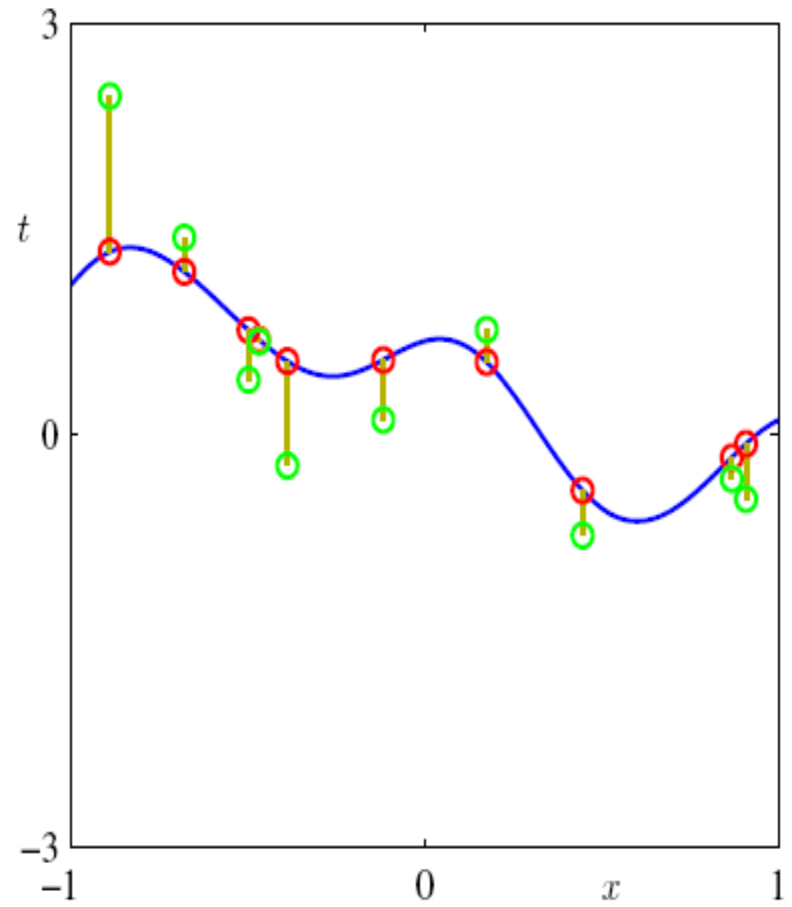


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

Gaussian processes for regression

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 $\{y_n\}$.



Gaussian processes for regression

- Goal of regression: predict $p(t_{N+1}|\mathbf{t}_N)$
- Using the joint distribution

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

- Partitioning the covariance

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

- \mathbf{K} has elements $k(\mathbf{x}_n, \mathbf{x}_{N+1})$ and c is given as

$$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

- $p(t_{N+1}|\mathbf{t})$ is Gaussian with mean and covariance given by

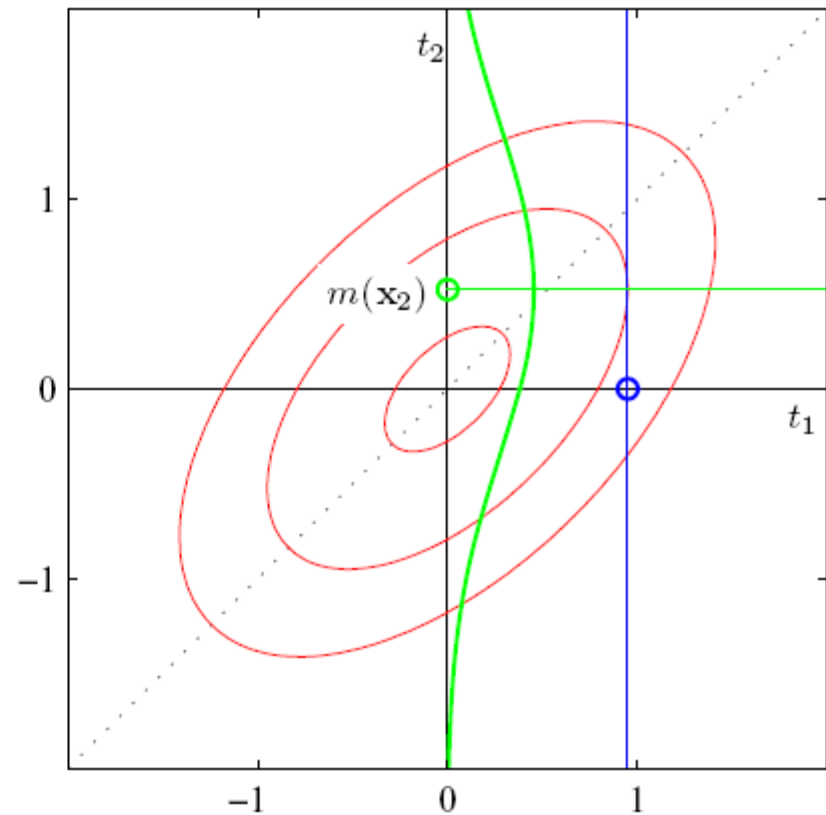
$$\begin{aligned} m(\mathbf{x}_{N+1}) &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \\ \sigma^2(\mathbf{x}_{N+1}) &= c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \end{aligned}$$

Gaussian processes for regression

- Note that the mean and covariance are dependent on the term \mathbf{k} which is dependent on the input \mathbf{x}_{N+1}
- Note also that the additional kernel matrix should be a valid kernel.
- The Gaussian process viewpoint is advantageous in that we can consider covariance functions that can be expressed in terms of an infinite number of basis functions.

Gaussian processes for regression

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.



Gaussian processes for regression

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.

