# Chapter 3  Linear Models for Regression

## Regression (supervised learning)

$x$ :  D-dimensional  input  vector

$t$ :  continuous  target  variable

Linear regression model : linear combination of non linear basis functions of the input variables with adjustable parameters.  E.g.  of basis function

$$y(x, w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$x \xrightarrow{\Phi} \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix}$$
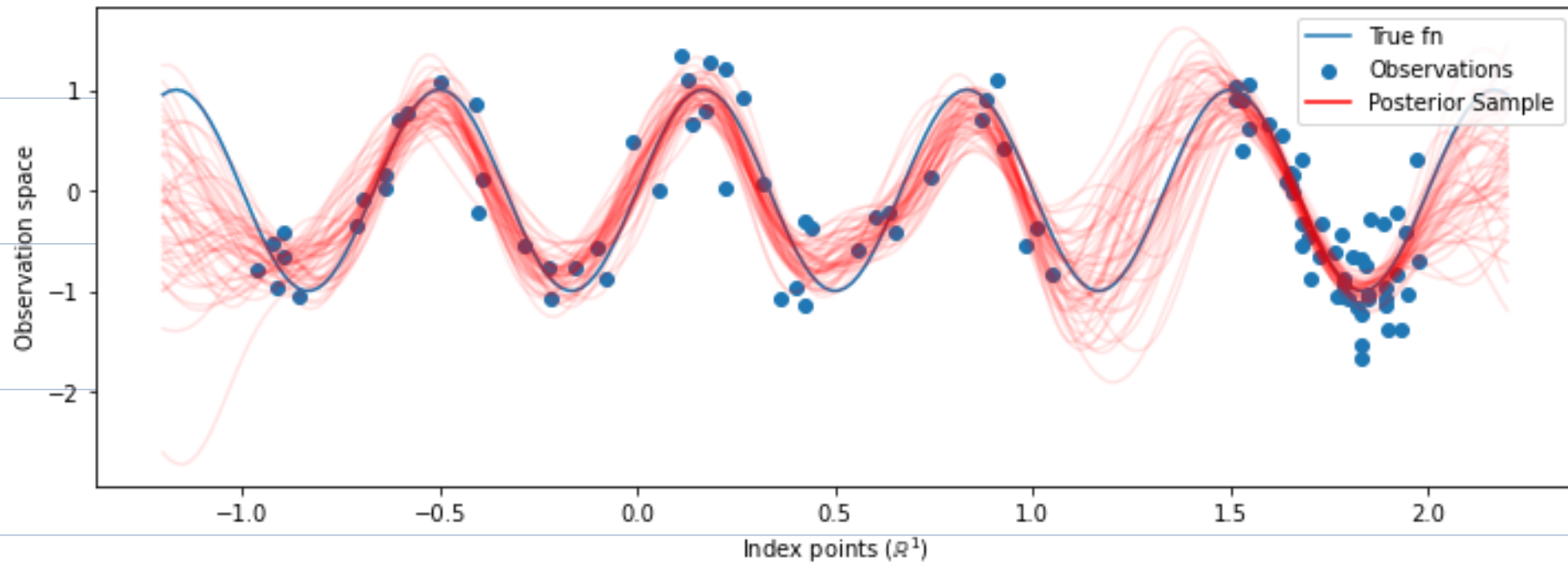
whose basis is $\{1, x, x^2, x^3\}$.

$N$ observations $\{x_n\}$ $n = 1, 2, \ldots N$, corresponding target $t_n$.

Goal: predict the value of $t$ for a new $x$.

Simplest approach: constructing appropriate function $y(x)$

General or probabilistic perspective: modeling the predictive distribution $p(t \mid x)$ (uncertainty about $t$ for each $x$)

# 3.1 Linear Basis Function Models

Linear combinations of fixed nonlinear functions of input $x$

$$y(x, w) := w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x)$$

where $\phi_j(x)$ are known as basis functions. $w_0$ allows for any fixed offset, called bias.

So it is convenient to define an additional dummy 'basis function' $\phi_0(x) = 1$ so that

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \Phi(x)$$

D-dim    M-dim
$$x \rightarrow \Phi(x)$$

where $w = (w_0, \ldots w_{M-1})^T$ and $\Phi(x) = (\underset{=1}{\underline{\phi_0}}, \phi_1(x), \ldots \phi_{M-1}(x))^T$

In view of pre-processing or feature extraction, the

feature can be expressed as $\{\phi_j(x)\}$

## Basis functions

In Chapter 1. there is a single input $x$ and the basis

functions take the form of powers of $x$ ( $\phi_j(x) = x^j$ )

One limitation of polynomial basis : global functions

# Gaussian basis functions

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

where $\mu_j$ govern the locations of the basis functions in input space and $s$ governs their spatial scale.

# Sigmoidal basis functions

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

where $\sigma$ is the logistic sigmoid $\sigma(a) = \frac{1}{1+\exp(-a)}$

Equivalently, we can use the 'tanh' function. Since $\tanh(a) = 2\sigma(2a) - 1$, general linear combination of sigmoid is equivalent to a general linear combination of tanh

Most of the discussion in this chapter is independent of the particular choice of basis.

### 3.1.1 Maximum likelihood and least squares

We have showed SSE could be motivated as the maximum likelihood solution under an assumed Gaussian noise model. As before, we assume that target $t$ is given by a deterministic function $y(x, w)$ with additive Gaussian noise

$$t = y(x, w) + \varepsilon$$

where $\varepsilon$ is a zero mean Gaussian random variable with precision $\beta$ (inverse of variance)

I.e.

$$p(t \mid x, w, \beta) = N(t \mid y(x, w), \beta^{-1}) \qquad (3.8)$$

In section 1.5.5, we showed that

$$E_t[t \mid x] = \int t \, p(t \mid x) \, dt \qquad (3.9)$$

is the optimal prediction

Consider inputs $\mathbb{X} = \{x_1, \ldots x_N\}$ with corresponding target

$t_1, \ldots t_N$. Let $\mathbb{t} := (t_1, \ldots, t_N)^T$. Assume $\mathbb{X}$ and $\mathbb{t}$ are

drawn independently from (3.8). Then the likelihood function

of $w$ and $\beta$ is in the form

$$p(\mathbb{t} \mid \mathbb{X}, w, \beta) = \prod_{n=1}^{N} \mathcal{N}(t_n \mid \underbrace{w^T \Phi(x_n)}_{y(x, w)}, \beta^{-1}) \qquad (3.10)$$

We will drop the explicit $\mathbb{X}$ from expressions.

$$\ln p(\mathbb{t} \mid w, \beta) = \sum_{n=1}^{N} \ln \mathcal{N}(t_n \mid w^T \Phi(x_n), \beta^{-1}) \qquad (3.11)$$

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(w)$$

where SSE is defined by

$$E_D(w) := \frac{1}{2} \sum_{n=1}^{N} \{ t_n - w^T \Phi(x_n) \}^2 \qquad (3.12)$$

실제 예측

Consider first the maximization of (3.11) w.r.t $w$.

Maximization of likelihood function under a conditional Gaussian

(3.10) for a linear model $\iff$ Minimizing $E_D(w)$

$$\nabla_w \ln p(t | w, \beta) = \beta \sum_{n=1}^{N} \{ t_n - w^T \Phi(x_n) \} \Phi(x_n)$$

Setting this gradient to zero gives

$$0 = \sum_{n=1}^{N} t_n \Phi(x_n) - \left( \sum_{n=1}^{N} \Phi(x_n) \Phi(x_n)^T \right) w$$

Solving for $w$ we obtain

$$\Phi(x_n) := \begin{pmatrix} \phi_1(x_n) \\ \vdots \\ \phi_{M-1}(x_n) \end{pmatrix}$$

$$w_{ML} = (\underset{M \times N \; N \times M}{\Phi^T \Phi})^{-1} \underset{M \times N}{\Phi^T} \underset{N \times 1}{t} \qquad O(M^3)$$

which are known as normal equations for least square problem.

Here $\Phi$ is an $N \times M$ matrix called design matrix whose

elements are given by $\Phi_{nj} = \underline{\phi_j(x_n)}$

I.e.

$$
\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & & \phi_{M-1}(x_2) \\ \vdots & \vdots & & \\ \phi_0(x_N) & & & \phi_{M-1}(x_N) \end{pmatrix} = \begin{pmatrix} \Phi(x_1)^T \\ \Phi(x_2)^T \\ \vdots \\ \Phi(x_N)^T \end{pmatrix}
$$

where $\quad \Phi(x) = (\phi_0(x), \cdots \phi_{M-1}(x))^T$.

The quantity

$$
\Phi^\dagger := (\Phi^T \Phi)^{-1} \Phi^T
$$

is known as the Moor - Penrose pseudo - inverse of $\Phi$.

If $\Phi$ is square and invertible, then $\Phi^+ = \Phi^{-1}$.

Let us see the role of bias parameter $w_0$

$w_0$

$$E_D(w) = \frac{1}{2} \sum_{n=1}^{N} \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(x_n) \right\}^2$$

SSE

Set $\frac{\partial E_D}{\partial w_0} = 0$ and solving for $w_0$. Then we obtain

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j$$

where we have defined $\quad \bar{t} := \frac{1}{N} \sum t_n \quad, \quad \bar{\phi}_j := \frac{1}{N} \sum \phi_j(x_n)$

Thus, $w_0$ is the difference between the averages of $t_n$ and the weighted sum of the averages of basis function values

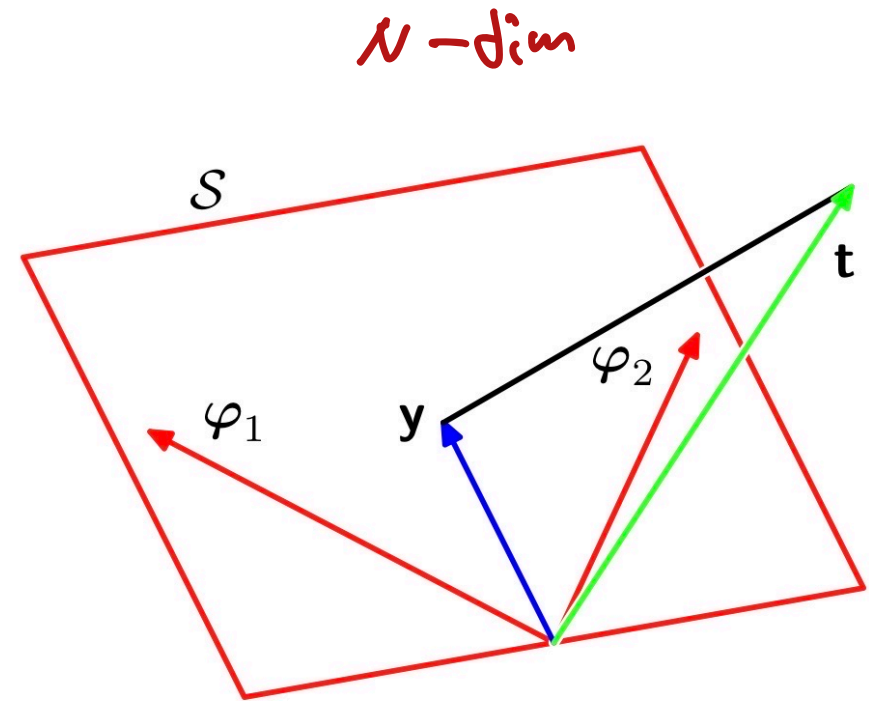After finding $w_{ML}$, we can maximize log likelihood (3.11) w.r.t noise precision parameter $\beta$,

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^{N} \underbrace{\left\{ t_n - w_{ML}^T \phi(x_n) \right\}^2}_{\text{residual variance of target}}$$

# 3.1.2 Geometry of Least squares

**Figure 3.2** Geometrical interpretation of the least-squares solution, in an $N$-dimensional space whose axes are the values of $t_1, \ldots, t_N$. The least-squares regression function is obtained by finding the orthogonal projection of the data vector $\mathbf{t}$ onto the subspace spanned by the basis functions $\phi_j(\mathbf{x})$ in which each basis function is viewed as a vector $\varphi_j$ of length $N$ with elements $\phi_j(\mathbf{x}_n)$.

$M > H$



$$\mathbf{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$$

$N$ - dimensional space whose axes are given by $t_n$.

For fixed $j$, basis function values $\phi_j(\mathbf{x}_n)$ ( $N$ data points) can be represented as a vector in the same space

$N-\dim$

Denote $\varphi_j$ by this vector, given by

$$\varphi_j := ( \phi_j(x_1), \phi_j(x_2), \dots , \phi_j(x_N))^T \quad ( j^{th} \text{ colum of } \Phi )$$

where $j = 0, 1, \dots, M-1$. <span style="color:red">N 차원 M 개의 vector</span>

Let $M < N$ and $S$ be the $M$-dim subspace spanned by $\varphi_j$

Define $y := ( \underbrace{Y(x_1, w)}_{= \sum\limits_{j=0}^{M-1} w_j \phi_j(x_i)}, Y(x_2, w), \dots Y(x_N, w))^T$. Because $y$

<span style="color:red">☆ $y = w_0 \varphi_0 + w_1 \varphi_1 + \dots w_{M-1} \varphi_{M-1}$ ☆</span>

is an arbitrary linear combination of $\varphi_j$, $y$ live anywhere

in the $M$-dimensional subspace $S$. <span style="color:red">y 의 모든 component는</span>

<span style="color:red">동일한 $w_0 \dots w_{M-1}$ 을 공유</span>

SSE (3.12) is equal (up to a factor $\frac{1}{2}$) to $\|\tilde{y} - \sharp\|^2$

( squared Euclidean distance)

Thus the least square solution $w$ corresponds to that

choice of $\tilde{y}$ lying in subspace $S$ and that is closest to $\sharp$

### 3.1.3 Sequential learning

A.k.a on-line algorithm

Applying the technique of stochastic gradient descent, also known as sequential gradient descent

If the error function $E = \sum_n E_n$, then after presentation of pattern $n$, SGD updates $w$ using

$$w^{(\tau+1)} := w^{(\tau)} - \eta \nabla E_n$$

where $\tau$: iteration number, $\eta$: learning rate parameter.

For the case of SSE (3.12),
$$E_n = \frac{1}{2}\left(t_n - w^T \Phi(x_n)\right)^2$$

$$w^{(\tau+1)} = w^{(\tau)} + \eta \underbrace{\left(t_n - w^{(\tau)T} \overset{\text{실제}}{\Phi_n}\right)}_{\text{예측}} \Phi_n$$

where $\overline{\Phi}_n := \Phi(x_n) = \left(\phi_0(x_n),\ \phi_1(x_n),\ \dots\ \phi_{M-1}(x_n)\right)^T$

# 3.1.4 Regularized least squares

To prevent over-fitting, we added regularization term so that

$$\overset{SSE}{E_D(w)} + \lambda E_w(w)$$

where $\lambda$ is the regularization coefficient.

Simple for of regularizer is as follow

$$E_w(w) := \frac{1}{2} w^T w$$

weight decay

parameter shrinkage

If we also consider SSE, then the total error function becomes

$$\frac{1}{2} \sum_{n=1}^{N} \{ t_n - w^T \Phi(x_n) \}^2 + \frac{\lambda}{2} w^T w \qquad (3.27)$$

Set the gradient of (3.27) w.r.t $w$ to zero and solve for $w$. Then we obtain the solution $w$

$$w = (\lambda I + \Phi^T \Phi)^{-1} \Phi \ t$$

More general regularizer is used as follows

$$\frac{1}{2}\sum_{n=1}^{N}\{t_n - w^T\Phi(x_n)\}^2 + \frac{\lambda}{2}\sum_{j=1}^{M}|w_j|^q \qquad (3.29)$$

where $q=2$ corresponds to the quadratic regularizer (3.27)

The case of $q=1$ is known as lasso.

Excercise 3.5 and Appendix E

Minimize (3.29) $\iff$ Minimize $E_D(w/)$ subject to the constraints

$$\sum_{j=1}^{M} |w_j|^q \le \eta$$



for some appropriate value of $\eta(\lambda)$

**Figure 3.4** Plot of the contours of the unregularized error function (blue) along with the constraint region (3.30) for the quadratic regularizer $q = 2$ on the left and the lasso regularizer $q = 1$ on the right, in which the optimum value for the parameter vector $\mathbf{w}$ is denoted by $\mathbf{w}^\star$. The lasso gives a sparse solution in which $w_1^\star = 0$.

## 3.1.5 Multiple outputs

$\boxed{K > 1}$ dimensional target vector $t = \begin{pmatrix} t_1 \\ \vdots \\ t_K \end{pmatrix}$

Our approach is to use the same basis functions to model

all of the components of target vector

$$ y(x, W) := W^{T}\Phi(x) = \begin{pmatrix} y_1 \\ \vdots \\ y_K \end{pmatrix} $$

$\nearrow$ M-dim

where $y$ is a $K$-dim vector, $W$ is an $M \times \boxed{K}$ matrix

of parameters and $\Phi(x)$ is an $M$-dim vector with elements

$\phi_j(x)$

Suppose conditional distribution of the target vector to be an _istropic_ Gaussian

$$p(t \mid x, W, \beta) = \mathcal{N}(t \mid W^T \phi(x), \beta^{-1} I)$$

single value (pointing to $\beta$)

Given $k$-dim $N$ observations $t_1, t_2, \ldots t_N$. we can combine these into $N \times k$ matrix $T$. (target)

Similarly combine the input vectors $x_1, x_2, \ldots x_N$ into $N \times D$ matrix $X$.

The log likelihood

$$\ln p(T \mid X, W, \beta) = \sum_{n=1}^{N} \ln \mathcal{N}(t_n \mid W^T \Phi(x_n), \beta^{-1} I)$$

$$= \frac{Nk}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{n=1}^{N} \| t_n - W^T \Phi(x_n) \|^2$$

Maximization solution for $W$ is given by

$$W_{ML} = \underbrace{(\underbrace{\Phi^T \Phi}_{M \times N \ N \times M})^{-1} \underbrace{\Phi^T}_{M \times N} \underbrace{T}_{N \times K}} \qquad M \times K \quad \text{matrix}$$

If we examine this result for each target variable $t_k$,

basis function values of observations

$$W_k = (\Phi^T \Phi)^{-1} \Phi^T t_k = \overbrace{\Phi^\dagger} \underline{t_k}$$

target of observations

where $t_k$ is an $N-$dim column vector ($k^{th}$ column of $T$)

Thus, the solutions decouples between the different targets.

From now on, we will consider single target variable $t$.

$$x \in \mathbb{R}^D \quad \text{input}, \quad \text{basis} \quad \text{function} \quad \Phi(\cdot) = \begin{pmatrix} 1 \\ \phi_1(x) \\ \vdots \\ \phi_{M-1}(x) \end{pmatrix}$$

target $t$, determine $w$ $\overset{M \times 1}{}$

$$y(x, w) = w^T \Phi(x) \approx t \qquad \underset{\text{1-dim Gaussian}}{( t \sim N(t \mid y(x,w), \beta )}$$

target $t\!\!t$ ( $K$ - dim ), determine $W$ $\overset{M \times K}{}$

$$y(x, W) = W^T \Phi(x) \approx t\!\!t \qquad \underset{\text{K-dim Gaussian}}{( t\!\!t \sim N(t\!\!t \mid y(x,w), \beta I )}$$

# 3.2 The Bias - Variance Decomposition <span style="color:red">( 모델 평가)</span>

## Frequentist view of model complexity

## Bias - Variance trade-off

When error function is SSE, the optimal prediction is given by

$$h(x) = \mathbb{E}_t[t \mid x] = \int t \, p(t \mid x) \, dt$$

<span style="color:red">SSE</span>  <span style="color:red">$L(x,t)$</span>

We showed in Section 1.5.5 that the expected squared loss can be written in the form

(3.37)

<span style="color:red">Prediction</span>  <span style="color:red">optimal solution</span>

$$\mathbb{E}[L] = \int \{ y(x) - h(x) \}^2 \, p(x) \, dx \; + \; \int \{ h(x) - t \}^2 \, p(x,t) \, dx \, dt$$

$$\mathbb{E}[L] = \int \{ \overset{\text{Prediction}}{y(x)} - \overset{\text{optimal solution}}{h(x)} \}^2 p(x) dx \;+\; \int \{ h(x) - t \}^2 p(x,t) dx \, dt$$

$$\overset{\downarrow}{y(x, w)}$$

The second term arises from the intrinsic noise and
is the minimum expected loss.

The first term depends on our choice of $y(x)$

Our goal is to seek $y(x)$ making the first term a minimum.

유한한 $N$개의 관측

$\Rightarrow$ $D \rightarrow W_{ML} \longrightarrow y(x, w_{ML})$
매개변수

$y(x ; D)$
$\parallel$

modeling $h(x)$
using $y(x, w)$

Bayesian : uncertainty is expressed through a posterior distribution over $w$

Frequentist : point estimate of $w$ based on $D$.

$N$ observations $D$ are independently drawn from $p(t, x)$

For a given $D$, we can obtain a prediction function $y(x; D)$

$y(x; D)$ and its squared error depend on $D$

The performance of learning algorithm is assessed by taking the average over ensemble of data sets

Consider the first term in (3.37)

$$\{ y(x;D) - h(x) \}^2$$

prediction by $D$ depended by ML algorithm

optimal solution

which depends on $D$.

$$\{ y(x;D) \pm \mathbb{E}_D [ y(x;D)] - h(x) \}^2$$

$$= \{ y(x;D) - \mathbb{E}_D [ y(x;D)] \}^2 + \{ \mathbb{E}_D[y(x;D)] - h(x) \}^2$$

$$+ 2\{ y(x;D) - \mathbb{E}_D[y(x;D)]\} \{ \mathbb{E}_D[y(x;D)] - h(x)\}$$

Take the expectation w.r.t $D$.

$$E_D[\{y(x;D) - h(x)\}^2] = \underbrace{\{E_D[y(x;D)] - h(x)\}^2}_{\text{bias}^2}$$

$$+ \underbrace{E_D[\{y(x;D) - E_D[y(x;D)]\}^2]}_{\text{variance}}$$

bias : extent to which the average prediction over all data sets differs from the desired regression function

variance : extent to which the solutions for indivitual data sets vary around thier average.

( sensitivity of $y(x;D)$ w.r.t the choice of $D$ )

We can obtain the following decomposition of expected squared loss

$$\text{expected squared loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where
$$(\text{bias})^2 = \int \left\{ \mathbb{E}_D[y(x;D)] - h(x) \right\}^2 p(x)\, dx$$

$$\text{variance} = \int \mathbb{E}_D\left[ \left\{ y(x;D) - \mathbb{E}_D[y(x;D)] \right\}^2 \right] p(x)\, dx$$

$$\text{noise} = \int \left\{ h(x) - t \right\}^2 p(x,t)\, dx\, dt$$

Our goal is to minimize the expected loss.

Trade-off between bias and variance

Flexible models having high variance and low bias

Rigid models having low variance and high bias



$\ln \lambda = 2.6$

low variance

$\ln \lambda = -0.31$

high variance

$\ln \lambda = -2.4$

**Figure 3.5** Illustration of the dependence of bias and variance on model complexity, governed by a regularization parameter $\lambda$, using the sinusoidal data set from Chapter 1. There are $L = 100$ data sets, each having $N = 25$ data points, and there are $24$ Gaussian basis functions in the model so that the total number of parameters is $M = 25$ including the bias parameter. The left column shows the result of fitting the model to the data sets for various values of $\ln \lambda$ (for clarity, only 20 of the 100 fits are shown). The right column shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

Examine the bias - variance trade-off quantitavely

L prediction models $y^{(\ell)}$, $\ell = 1, \dots L$

$$D_1 \dots D_L$$

The average prediction

$$\bar{y}(x) := \frac{1}{L} \sum_{\ell=1}^{L} y^{(\ell)}(x)$$

and integrated squared bias and integrated variance

$$(bias)^2 = \frac{1}{N} \sum_{n=1}^{N} \left\{ \bar{y}(x_n) - h(x_n) \right\}^2$$

$$\left( \begin{array}{l} \text{approximated} \\ \text{by sum of } x_n \end{array} \right)$$

$$variance = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{L} \sum_{\ell=1}^{L} \left\{ y^{(\ell)}(x_n) - \bar{y}(x_n) \right\}^2$$

## 3.3 Bayesian Linear regression

Bayesian linear regression can avoid the overfitting problem of maximum likelihood and lead to automatic methods of determining model complexity.

$$t \sim N(t \mid y(x, w), \beta^{-1})$$

$x^*$ → predict $t^*$ ↗ parameter distribution

### 3.3.1 Parameter distribution

Consider the prior probability distribution over $w$

Noise precision parameter $\beta$ is assumed to be known

First, we noted likelihood $p(t \mid w)$ is the exponential

of quadratic of w

So the corresponding conjugate prior is given by

$$p(w) := \mathcal{N}(w \mid m_0, S_0)$$

where mean $m_0$, covariance $S_0$.

$$\prod_{n=1}^{N} \mathcal{N}(t_n \mid w^T \Phi(x_n), \beta^{-1})$$

Thus the posterior distribution in the form (see (2.116))

$$p(w \mid t) = \mathcal{N}(w \mid m_N, S_N)$$

where

$$m_N = S_N (S_0^{-1} m_0 + \beta \Phi^T t)$$

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

Since the posterior is Gaussian (unimodal), its mode = mean

$$\Rightarrow \quad w_{MAP} = m_N$$

If we consider $S_0 := \alpha^{-1} I$ and $\alpha \to 0$ i.e. infinitely

broad prior, then the mean $m_N$ reduces to $w_{ML}$

Similarly, if $N = 0$ (without observation), posterior = prior

For simplicity, consider a zero-mean isotropic Gaussian with single precision parameter $\alpha$ as a prior distribution

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

So the corresponding posterior

$$p(w|t) = \mathcal{N}(w|m_N, S_N)$$

where

$$m_N = \beta S_N \Phi^T t$$

$$S_N^{-1} = \alpha I + \beta \Phi^T \Phi$$

Log of posterior is the sum of log of likelihood and log of prior

$$\ln p(w \mid t) = -\frac{\beta}{2} \sum_{n=1}^{N} \{t_n - w^T \Phi(x_n)\}^2 - \frac{\alpha}{2} w^T w + \text{constant}$$

posterior

Its MAP solution w.r.t $w$ is equivalent to minimization of SSE with additional quadratic regularization term $\lambda = \alpha/\beta$.

Linear basis function

- input: $x$

- target: $t$

- $y(x, w) := w_0 + w_1 x$

- Observed data $\circ$

generated by $-0.3 + 0.5 x$

with std $0.2$

function of
$w$

$p(t \mid x, w)$

prior
$\mathcal{N}\left(\begin{smallmatrix} 0 \\ 0 \end{smallmatrix}, \begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix}\right)$
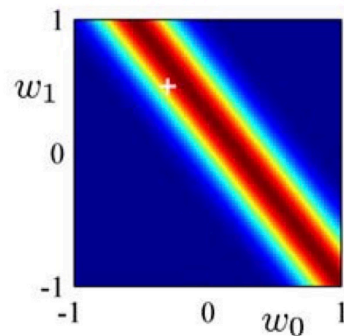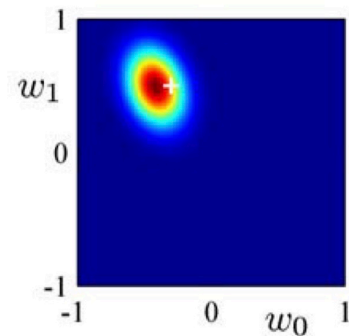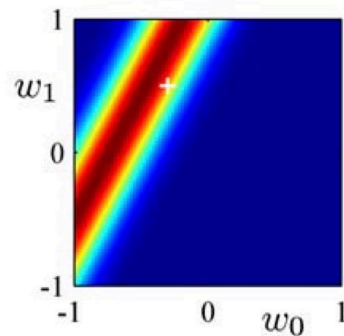
$(1, 0.1)$



**Figure 3.7** Illustration of sequential Bayesian learning for a simple linear model of the form $y(x, \mathbf{w}) = w_0 + w_1 x$. A detailed description of this figure is given in the text.

Generalized the Gaussian prior

$$p(w \mid \alpha) := \left[ \frac{q}{2} \left( \frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^{M} \exp \left( -\frac{q}{2} \sum_{j=0}^{M-1} |w_j|^q \right)$$

in which $q=2$ corresponds to the Gaussian.

If $q=2$, MAP solution of $w$ is the minimization solution

of $(3.29)$ which is SSE + regularization term

If $q \neq 2$, it is not true. (mode of posterior $\neq$ mean)

## 3.3.2 Predictive distribution

In practice, we are interested in making predictions of t
for a new $x$ ( not the value of $w$ )

Predictive distribution of t



$$p(t \mid x, \#, \alpha, \beta) = \int p(t \mid x, w, \beta) \, p(w \mid \#, \alpha, \beta) \, dw$$

<span style="color:red">↑ new input</span>     <span style="color:red">w/ on chặt posterior</span>

where $\#$ is the vector of training target values.

$\alpha$ is from prior assumption, $\beta$ is Gaussian noise of t

$$p(w \mid \alpha) = \mathcal{N}(w \mid 0, \alpha^{-1} I) \qquad\qquad p(t \mid x, w, \beta) = \mathcal{N}(t \mid y(x, w), \beta^{-1})$$

The predictive distribution takes the form

training data

$ \in R$

$$p(t \mid x, \pmb{t}, \alpha, \beta) = N(t \mid m_N^T \Phi(x), \sigma_N^2(x))$$

(
new input

where

$$\sigma_N^2(x) = \frac{1}{\beta} + \Phi(x)^T S_N \Phi(x) \qquad\qquad (3.59)$$

data noise

uncertainty
of w

By [ Qazaz et al., 1997]

$$\sigma_{N+1}^2(x) \leq \sigma_N^2(x)$$

If $N \to \infty$, then the second term in (3.59) $\to 1/\beta$
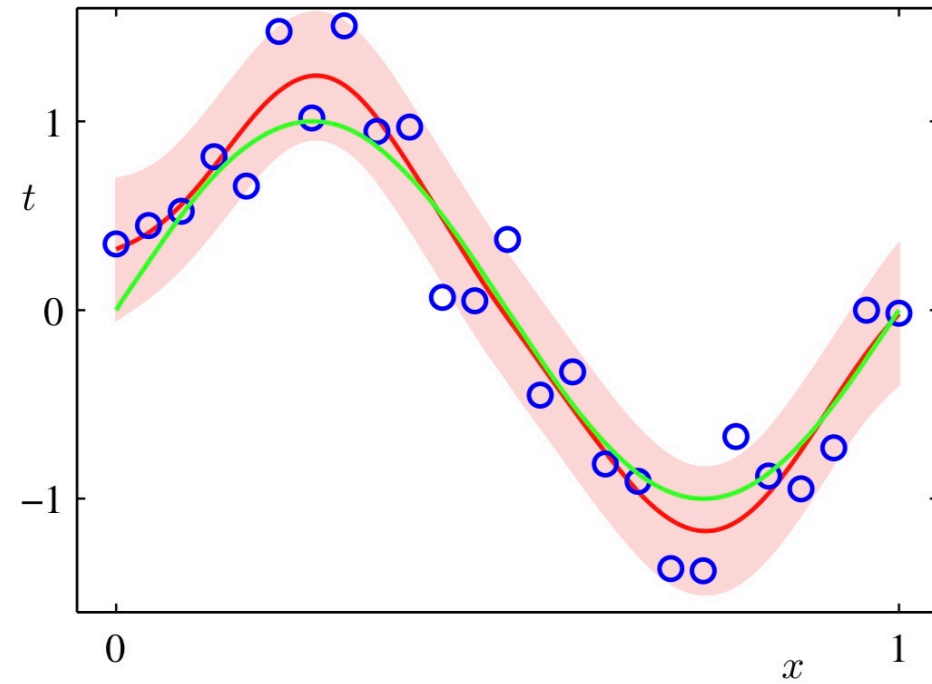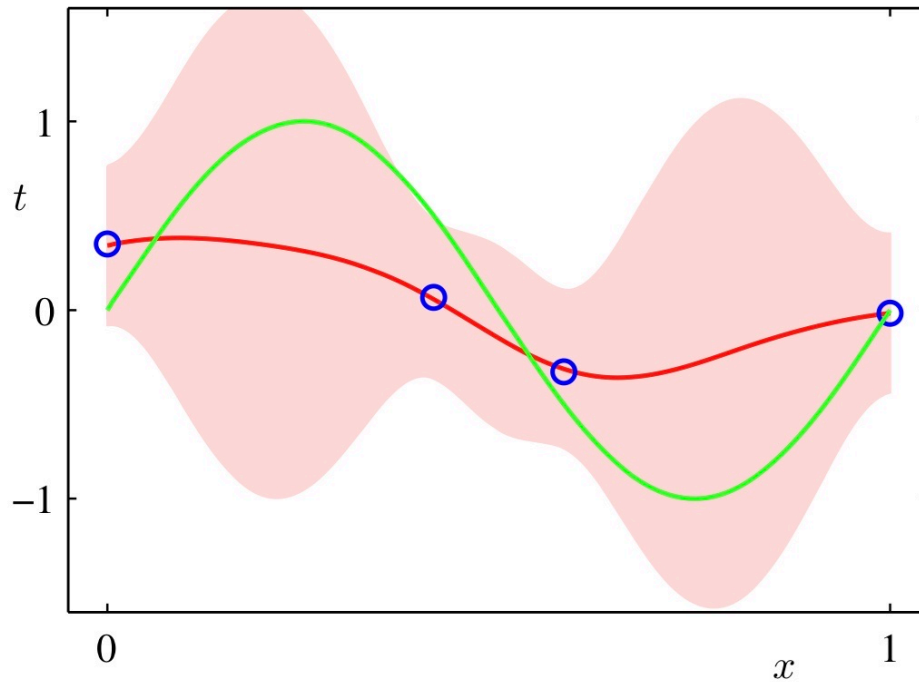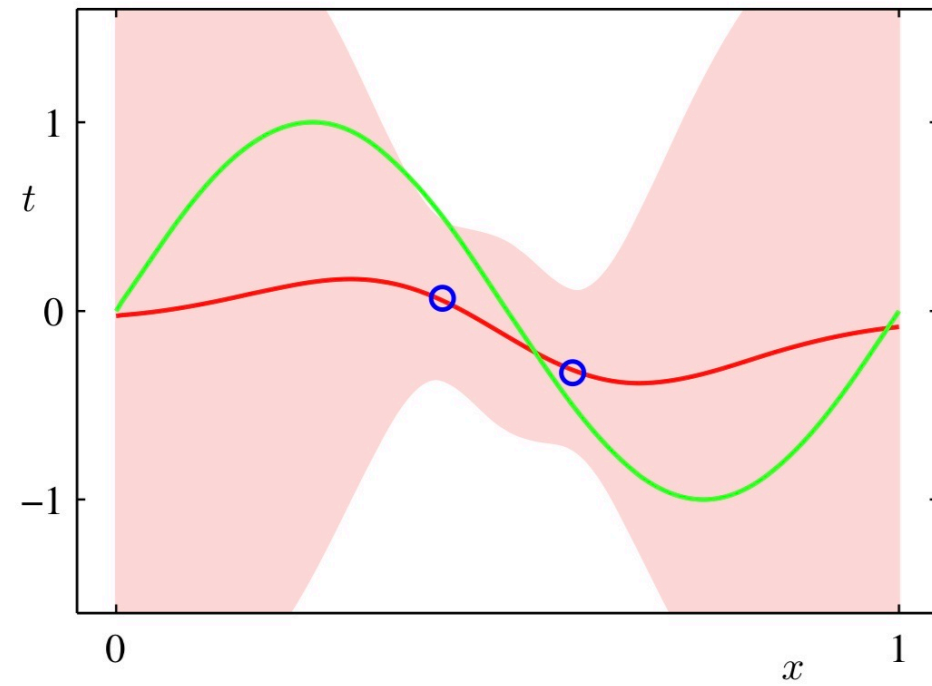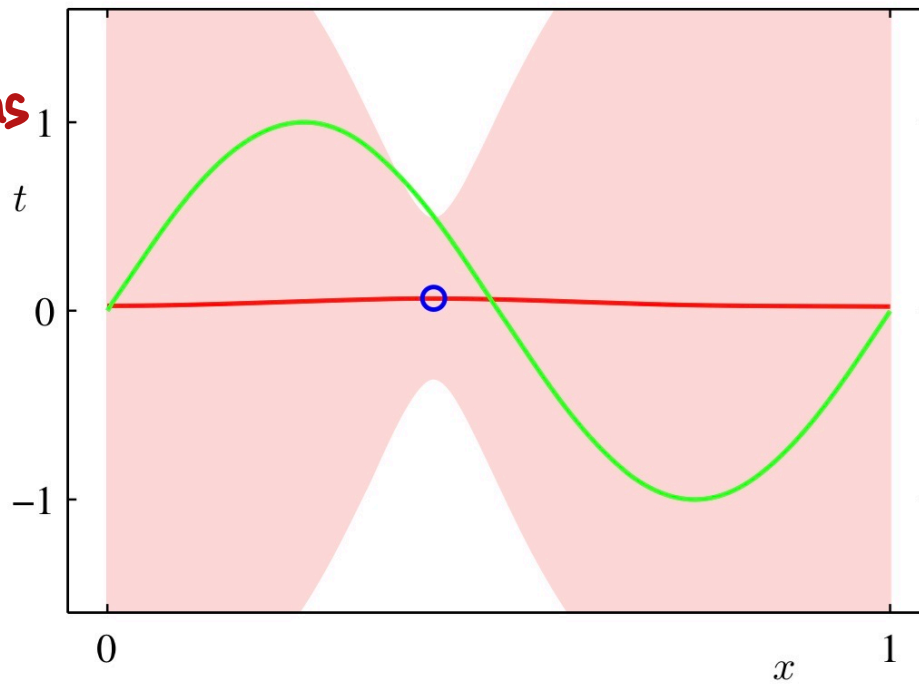
Uncertainty in predictions governed by (3.59)

**Figure 3.8** Examples of the predictive distribution (3.58) for a model consisting of 9 Gaussian basis functions of the form (3.4) using the synthetic sinusoidal data set of Section 1.1. See the text for a detailed discussion.

Sampling from posterior and plotting corresponding model
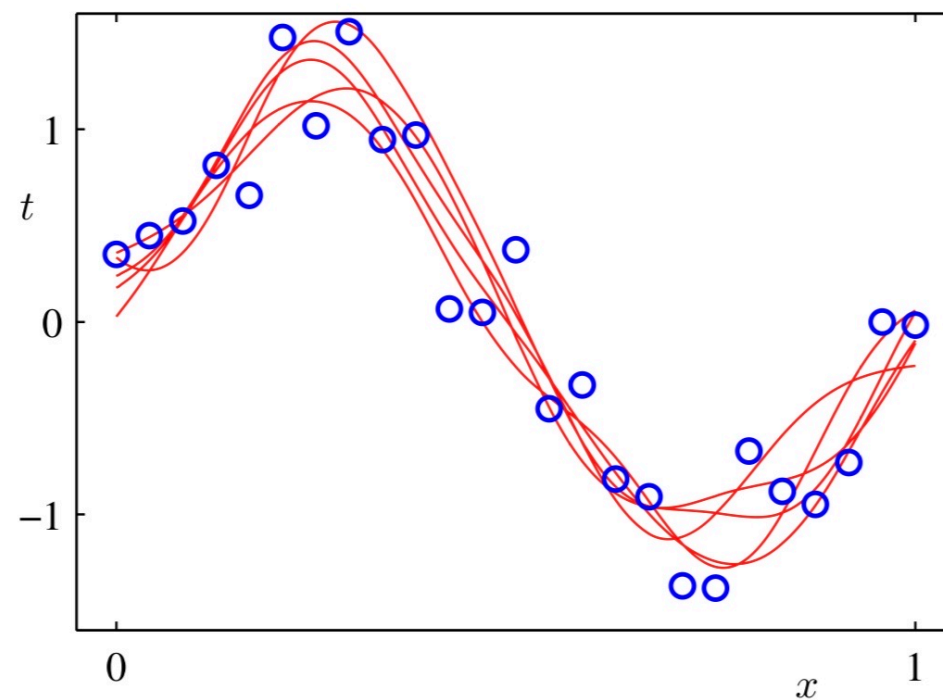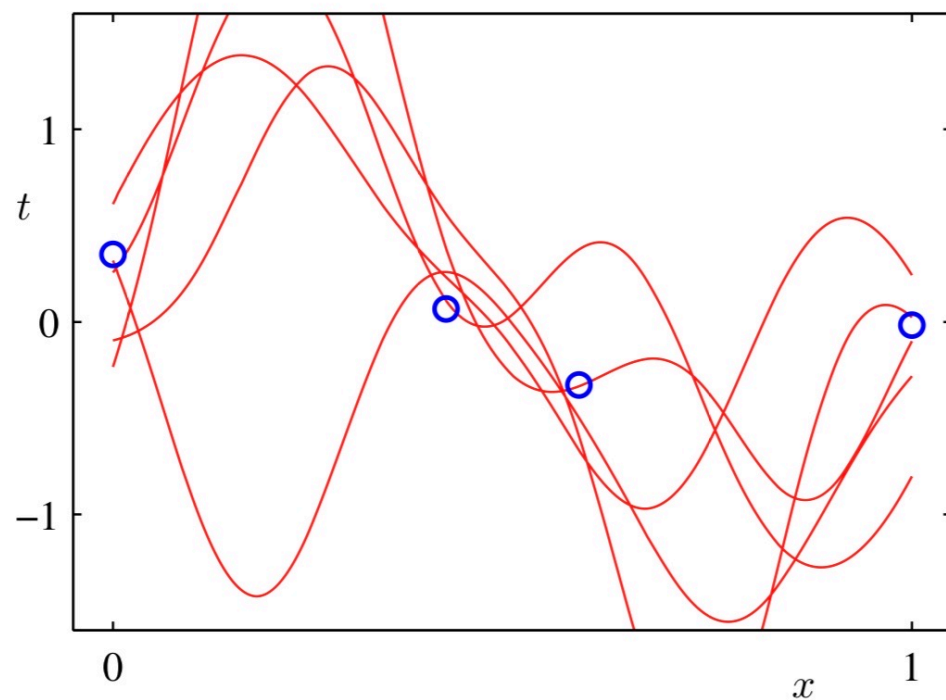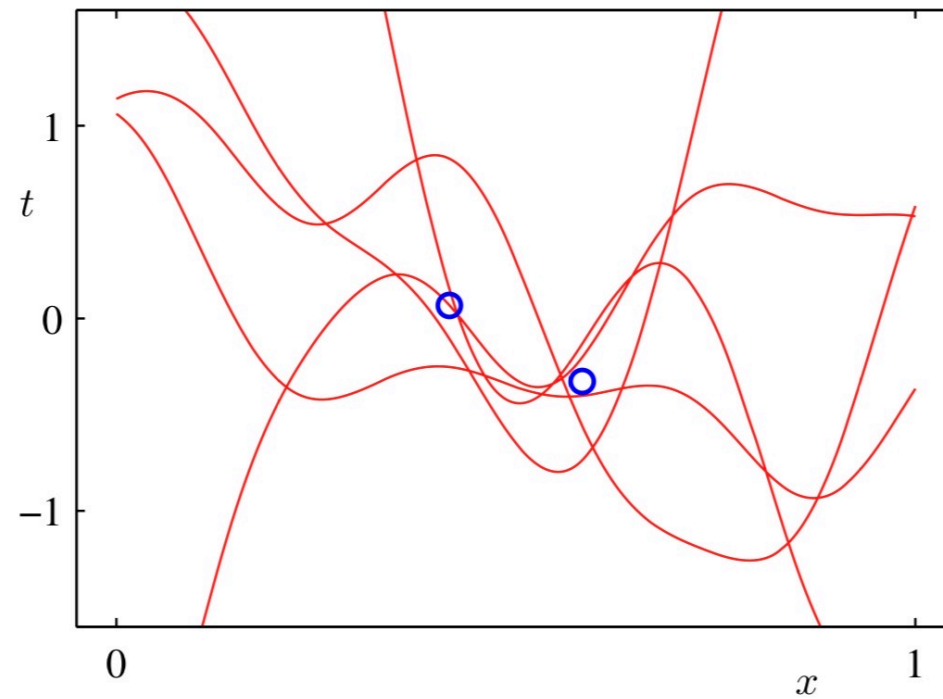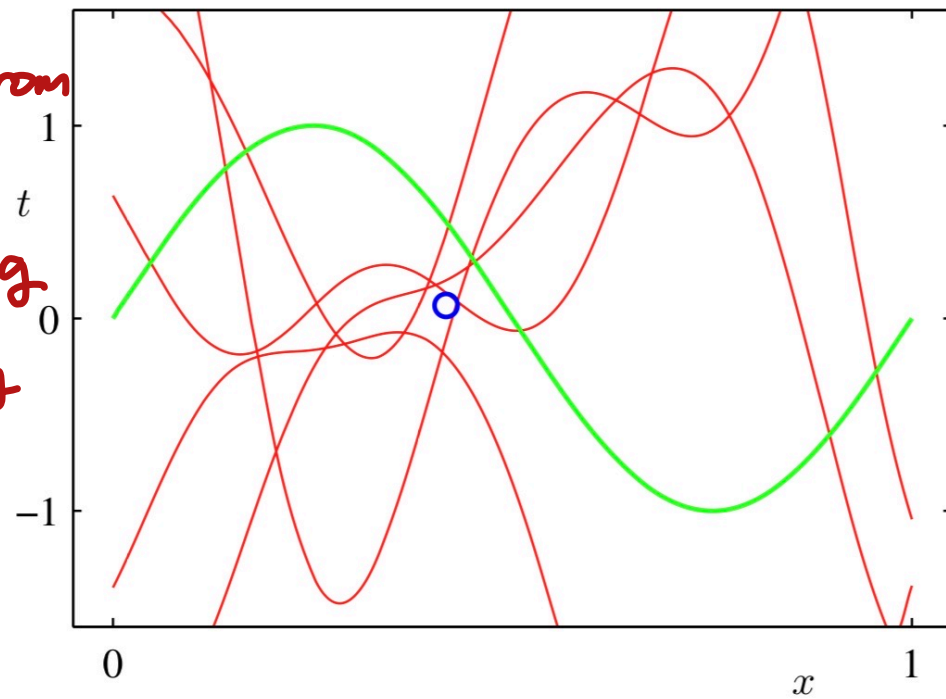
**Figure 3.9** Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distributions over $\mathbf{w}$ corresponding to the plots in Figure 3.8.

# Remark

- We have used Gaussian basis function (localized)

- If $x$ is away from the basis function centers, then the contribution from the second term in (3.59) goes to 0

  i.e. left the noise $\beta^{-1}$.

### 3.3.3 Equivalent kernel (kernel method)

Substitute (3.53) into (3.3) (expected prediction)

$$y(x, m_N) = m_N^T \Phi(x) = \beta \, \underset{M \times 1}{\Phi(x)}^T S_N \underset{M \times N}{\Phi^T} t = \sum_{n=1}^{N} \beta \, \Phi(x)^T S_N \Phi(x_n) \, t_n$$

where $\Phi(x) := (\phi_0(x), \dots \phi_{M-1}(x))^T$, $\qquad S_N^{-1} = S_0^{-1} + \beta \, \Phi^T \Phi$ and

$$\Phi = \begin{pmatrix} \phi_0(x_1) & & \phi_{M-1}(x_1) \\ \vdots & \dots & \\ \phi_0(x_N) & & \phi_{M-1}(x_N) \end{pmatrix} \quad \text{design matrix}$$

Thus, $y(x, m_N)$ is the linear combination of the training set target variables $t_n$.

$$\Rightarrow \quad y(x, m_N) = \sum_{n=1}^{N} k(x, x_n) t_n$$

where the function

$$k(x, x') := \beta \, \Phi(x)^T S_N \, \Phi(x')$$

is known as smoother matrix or equivalent kernel

Linear smoother : regression function makes predictions by taking linear combinations of training target values

This kernel depends on $x_n$ because of $S_N$

Consider the covariance between $y(x)$ and $y(x')$

$$\text{cov}[y(x), y(x')] = \text{cov}[w^T \Phi(x), w^T \Phi(x')]$$

$$= \Phi(x)^T S_N \Phi(x') = \beta^{-1} k(x, x')$$

$y(x) = $ mean of $\mathcal{N}(t \mid m_N^T \Phi(x), \beta^{-1} + \Phi(x)^T S_N \Phi(x))$

$y(x) = w^T \Phi(x)$    where    $w = \mathcal{N}(w \mid m_N, S_N)$

$$\text{cov}[w^T \Phi(x), w^T \Phi(x')] = \underbrace{\mathbb{E}[\Phi(x)^T w w^T \Phi(x')] - \Phi(x)^T m_N m_N^T \Phi(x')}$$

$$= \Phi(x)^T \underbrace{\mathbb{E}[w w^T]} \Phi(x')$$

$$= \text{cov}[w] + \mathbb{E}[w]\mathbb{E}[w]^T$$

$$= S_N + m_N m_N^T$$

For regression, we introduced a set of basis functions so equivalent kernel was implicitly determined.

But we can define a localized kernel directly and use this to make predictions

The equivalent kernel (3.62) can be expressed in the form an inner product w.r.t $\underline{\Psi}(x)$ of nonlinear functions

$$k(x, z) := \underline{\Psi}(x)^T \underline{\Psi}(z)$$

where $\underline{\Psi}(x) := \beta^{1/2} S_N^{1/2} \underline{\Phi}(x).$

## 3.4 Bayesian Model Comparison

The problem of model selection from a Bayesian perspective.

The over-fitting associated with Maximum likelihood can be avoided by marginalizing over model parameters.

The Bayesian view of model comparison involves the use of probabilities to represent uncertainty in the choice of model.

Compare a set of L models $\{ M_i \}_{i=1,...L}$

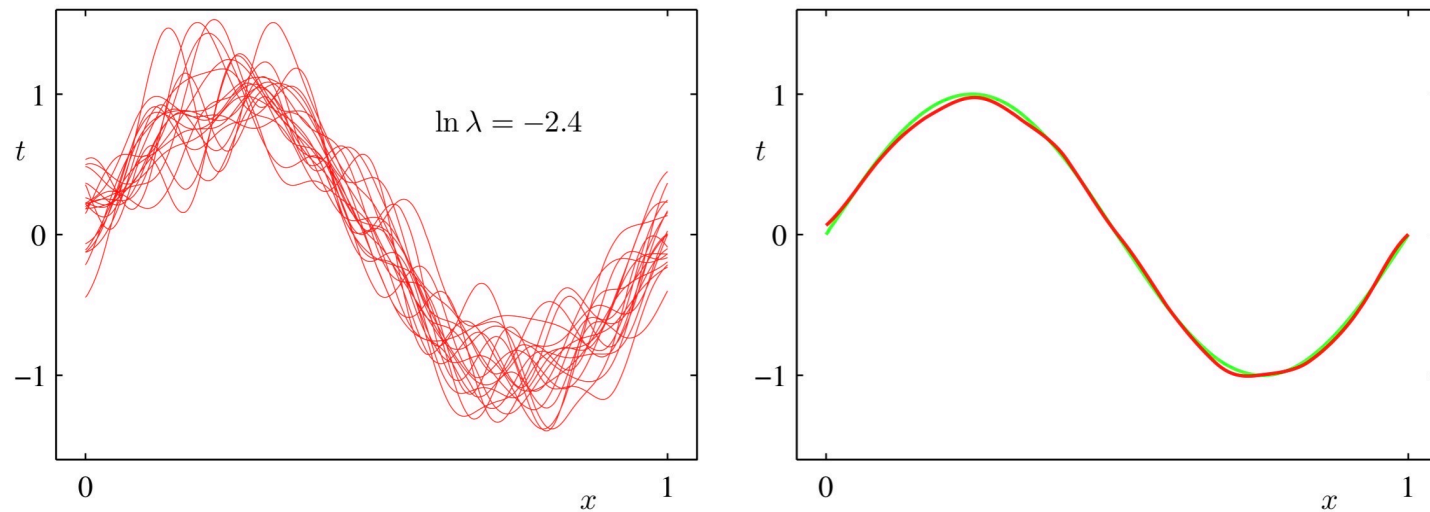Model refers to a probability distribution over the observed data D



**Figure 3.5** Illustration of the dependence of bias and variance on model complexity, governed by a regularization parameter $\lambda$, using the sinusoidal data set from Chapter 1. There are $L = 100$ data sets, each having $N = 25$ data points, and there are $24$ Gaussian basis functions in the model so that the total number of parameters is $M = 25$ including the bias parameter. The left column shows the result of fitting the model to the data sets for various values of $\ln \lambda$ (for clarity, only 20 of the 100 fits are shown). The right column shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

Suppose that the data is generated from one of models but we are uncertain which one

This uncertainty is expressed through $p(M_i)$

Given a training data set $D$, we want to evaluate

$$p(M_i \mid D) \propto \underbrace{p(M_i)}_{\text{prior}} \underbrace{p(D \mid M_i)}_{\text{model evidence}}$$

$$\theta$$
$$p(D \mid \theta)$$

The prior can express a preference for different models. But for simplicity assume that all models have the same prior.

$p(D \mid M_i)$ model evidence ( marginal likelihood ) expresses the preference shown by the data for different models.
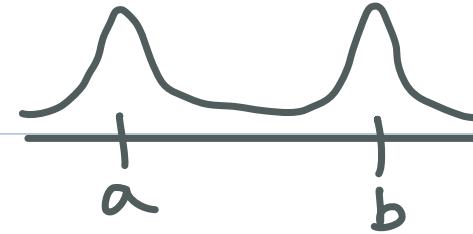
( likelihood function over the model space in which the parameters have been marginalized out )

The predictive distribution (mixture distribution)

$$p(t \mid x, D) = \sum_{\tilde{\lambda}=1}^{L} p(t \mid x, M_{\tilde{\lambda}}, D) \, p(M_{\tilde{\lambda}} \mid D)$$

posterior over models

Average of the predictive distributions $p(t \mid x, M_{\tilde{\lambda}}, D)$ of individual

models weighted by the posterior probabilities $p(M_{\tilde{\lambda}} \mid D)$

For example, two models

$M_1$

$M_2$



Model selection: use the single most probable model alone.

Consider model $M_i$ governed by the parameter $w$. The model evidence is given by

$$p(D \mid M_i) = \int p(D \mid w, M_i) \, p(w \mid M_i) \, dw$$

The model evidence (marginal likelihood) $p(D \mid M_i)$ can be viewed as the probability of generating the data set $D$ from a model whose parameters are sampled at random from the prior

Note that

$$P(w \mid D, M_i) = \frac{P(D \mid w, M_i) \, P(w \mid M_i)}{P(D \mid M_i)}$$

The model evidence is the normalization term appearing in the denominator in Bayes Theorem when evaluating the posterior w/

Consider a single parameter $w$. The posterior distribution over $w$ is proportional to $p(D|w)p(w)$.

For simplicity, assume the posterior distribution is sharply peaked around the most probable value $w_{MAP}$ with width $\Delta w_{posterior}$

**Figure 3.12** We can obtain a rough approximation to the model evidence if we assume that the posterior distribution over parameters is sharply peaked around its mode $w_{MAP}$.



and the prior is flat with width $\Delta w_{prior}$ so that

$$p(w) = 1/\Delta w_{prior}.$$

Thus we have a simple approximation to the integral over $w$

$$p(D) = \int p(D|w)\, p(w)\, dw \simeq p(D|w_{MAP}) \frac{\Delta w_{posterior}}{\Delta w_{prior}}$$

So

$$\ln p(D) \simeq \underbrace{\ln p(D|w_{MAP})}_{\text{fit to data given by the most probable parameter}} + \underbrace{\ln \left( \frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)}_{\text{penalty of the model complexity}}$$

<span style="color:red">fit to data given by the most probable parameter</span>

<span style="color:red">penalty of the model complexity</span>

$\Delta w_{posterior} < \Delta w_{prior}$, then the second term is negative.

So it increases in magnitude as the ratio $\Delta w_{post}/\Delta w_{prior}$ gets

smaller. If the parameters are finely tuned to the data in

posterior, then the penalty term is large.

For a model with M parameters, assume all parameters have the same ratio of $\Delta w_{posterior} / \Delta w_{prior}$, then we obtain a similar approximation as follows

$$\ln P(D) \simeq \ln P(D | w_{MAP}) + M \ln \left( \frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)$$

Thus, the size of the complexity penalty increases linearly

# 3.5 The evidence Approximation

Fully Bayesian treatment of linear basis function model

- Introduce prior distributions over hyperparameters $\alpha$ and $\beta$

- Make predictions by marginalizing w.r.t these hyperparameters and parameters w/

- But the complete marginalization over all of these variables $\alpha, \beta$ and w/ is analytically intractable.

Discuss an approximation in which we set the hyperparameters to specific values determined by maximizing the 'marginal likelihood function' obtained by first integrating over $w$.

If we introduce hyperprior over $\alpha$ and $\beta$, the predictive distribution is given by

$$p(t \mid \mathbb{t}) = \iiint p(t \mid w, \beta) \; p(w \mid \mathbb{t}, \alpha, \beta) \; p(\alpha, \beta \mid \mathbb{t}) \, dw \, d\alpha \, d\beta$$

(3.49)
posterior over $w$

model assumption
(3.8)

hyper posterior

Here we omitted the dependence on input $\mathbb{x}$.

If posterior $p(\alpha, \beta \mid t)$ is sharply peaked around $\hat{\alpha}$ and $\hat{\beta}$, then

$$p(t \mid t) \simeq p(t \mid t, \hat{\alpha}, \hat{\beta}) = \int p(t \mid w, \hat{\beta}) \, p(w \mid t, \hat{\alpha}, \hat{\beta}) \, dw$$

From Bayes Theorem, the posterior distribution for $\alpha, \beta$

$$p(\alpha, \beta \mid t) \propto p(t \mid \alpha, \beta) \, p(\alpha, \beta)$$

So if prior is relatively flat, the values $\hat{\alpha}$ and $\hat{\beta}$ are obtained by maximizing the marginal likelihood function $p(t \mid \alpha, \beta)$

Here, we evaluate the marginal likelihood for the linear basis model and then finding its maxima.

So this will allow us to determine values for hyperparameters from the training data alone. ($\alpha/\beta \simeq$ regularization parameter)

Two approaches of maximization of the log evidence

— Evaluate the evidence function analytically and then set its derivative equal to 0 to obtain re-estimation for $\alpha, \beta$

— Use the technique called expectation maximization algorithm in Section 9.3.4.

### 3.5.1 Evaluation of the evidence function

The marginal likelihood function $p(t \mid \alpha, \beta)$ is obtained by integrating over $w$

$$p(t \mid \alpha, \beta) = \int p(t \mid w, \beta) \, p(w \mid \alpha) \, dw$$

By the result (2.115) for the conditional distribution in a linear – Gaussian model, we can evaluate this integral.

From, (3.11), (3.12) and (3.52), we can write the evidence function in the form (Excercise 3.17)

$$p(\#|\alpha,\beta) = \left(\frac{\beta}{2\pi}\right)^{N/2}\left(\frac{\alpha}{2\pi}\right)^{M/2}\int \exp\{-E(w)\}\,dw \qquad (3.78)$$

where M is the dimensionality of w and

$$E(w) = \beta E_D(w) + \alpha E_w(w) \qquad (3.79)$$

$$= \frac{\beta}{2}\|\# - \Phi w\|^2 + \frac{\alpha}{2}w^T w$$

<span style="color:red">N    (N×M)(M×1)</span>

$$w \sim N(m_N, S_N^{-1})$$

Furthermore,

$$E(w) = E(m_N) + \frac{1}{2}(w - m_N)^T A (w - m_N)$$

where we have introduced

$$A = \alpha I + \beta \Phi^T \Phi$$

together with

$$E(m_N) = \frac{\beta}{2} \| t - \Phi m_N \|^2 + \frac{\alpha}{2} m_N^T m_N$$

A is the matrix of second derivatives of error function and a.k.a Hessian matrix

$$A = \nabla \nabla E(w).$$

Here we have also defined $m_N$ given by

$$m_N = \beta A^{-1} \Phi^T t. \qquad (3.84)$$

$$S_N^{-1} = \alpha I + \beta \Phi^T \Phi$$

Using (3.54), we see $A = S_N^{-1}$, hence (3.84) = (3.53)

Back to the integration (3.78)

$$\int \exp(-E(w)) \, dw$$

$$= \exp\{-E(m_N)\} \int \exp\left\{-\frac{1}{2}(w - m_N)^T A (w - m_N)\right\} \, dw$$

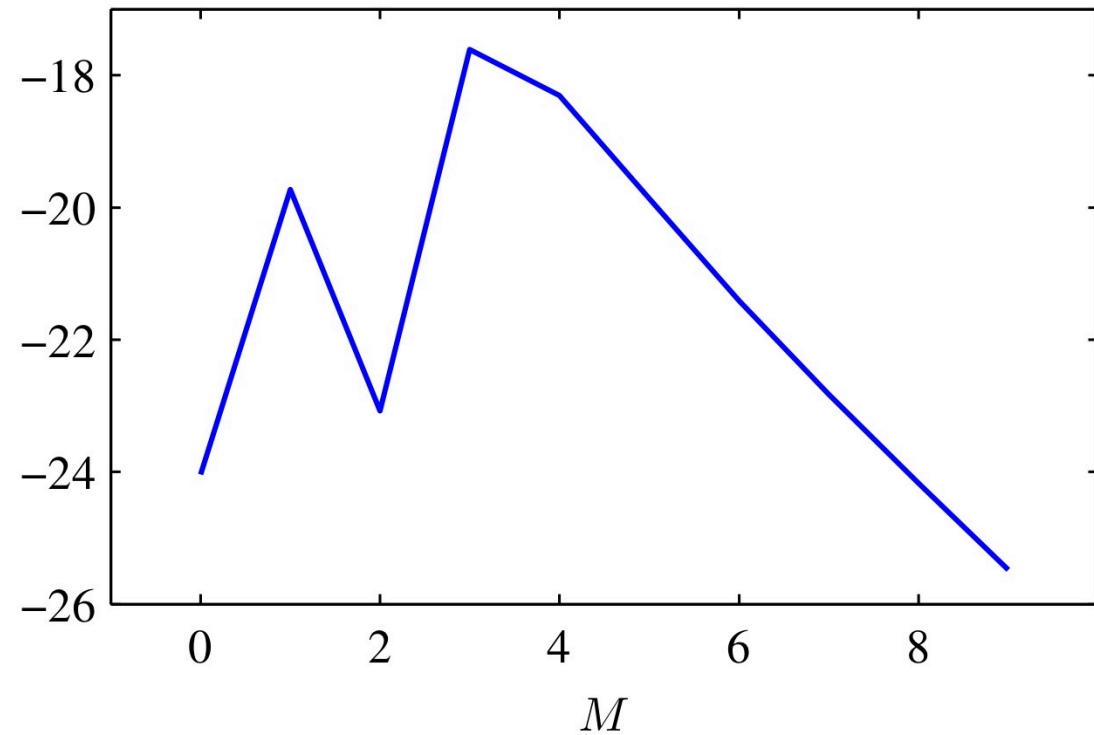$$= \exp\{-E(m_N)\} (2\pi)^{M/2} |A|^{-1/2}$$

Using (3.78), we can write the log of the marginal likelihood in the form

$$\ln p(t \mid \alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(m_N) - \frac{1}{2} \ln |A| - \frac{N}{2} \ln(2\pi)$$

**Figure 3.14** Plot of the log model evidence versus the order $M$, for the polynomial regression model, showing that the evidence favours the model with $M = 3$.



$$\ln p(\# \mid \alpha, \beta)$$

$$\alpha = 5 \cdot 10^{-3}$$

## Remark

- The underlying sinusoidal function is an odd function

- In $M = 3$ case, we obtain a significant improvement in data fit.

# 3.5.2 Maximizing the evidence function

Consider the maximization of $p(\boldsymbol{t}|\alpha,\beta)$ w.r.t $\alpha$.

This can be done by first defining the following eigenvector equation

$$(\beta \Phi^T \Phi) u_i = \eta_i u_i$$

Since $A = \alpha I + \beta \Phi^T \Phi$, $A$ has eigenvalues $\alpha + \eta_i$

Now consider the partial derivative of $\ln |A|$ w.r.t $\alpha$.

We have

$$\frac{\partial}{\partial \alpha} \ln |A| = \frac{\partial}{\partial \alpha} \ln \prod_{\lambda} (\eta_{\lambda} + \alpha) = \frac{\partial}{\partial \alpha} \sum_{\lambda} \ln (\eta_{\lambda} + \alpha) = \sum_{\lambda} \frac{1}{\eta_{\lambda} + \alpha}$$

Thus,

$$0 = \frac{M}{2\alpha} - \frac{1}{2} m_{N}^{T} m_{N} - \frac{1}{2} \sum_{\lambda} \frac{1}{\eta_{\lambda} + \alpha}$$

Multiplying by $2\alpha$ and rearranging, we obtain

$$\alpha \, m_{N}^{T} m_{N} = M - \alpha \sum_{\lambda} \frac{1}{\eta_{\lambda} + \alpha} =: \gamma$$

Since there are $M$ terms in the sum over $i$,

(3.91)
$$\gamma = \sum_i \frac{\eta_i}{\alpha + \eta_i} \qquad (\text{depends on } \alpha)$$

So the following $\alpha$ maximizes the marginal likelihood

$$\alpha = \frac{\gamma}{m_N^T \, m_N} \qquad (3.92)$$

Note that $\gamma$ depends on $\alpha$ and the mode $m_N$ of the posterior distribution depends on the choice of $\alpha$.

Thus, this solution is implicit and is adopted an iterative procedure

Make an initial choice of $\alpha$ and use this to find $m_N$ (3.53) and evaluate $\gamma$ (3.91)

Using (3.92), re-estimate $\alpha$ and the process repeat until convergence.

Note that because the matrix $\Phi^T \Phi$ is fixed, we can compute its eigenvalues once at the start

The value of $\alpha$ has been determined purely by training data.

Similarly, maximize the log marginal likelihood (3.86) w.r.t $\beta$.
Note that the eigenvalues $\lambda_i$ are proportional to $\beta$
and hence $d\lambda_i/d\beta = \lambda_i/\beta$ giving

$$\frac{\partial}{\partial\beta} \ln|A| = \frac{\partial}{\partial\beta} \sum_i \ln(\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta}$$

So the stationary point of the marginal likelihood satisfies

$$0 = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^{N} \{ t_n - m_N^T \phi(x_n) \}^2 - \frac{\gamma}{2\beta}$$

and rearranging we obtain

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^{N} \{ t_n - m_N^T \Phi(x_n) \}^2$$

Again, this is an implicit solution for $\beta$. So choose an initial value for $\beta$ and calculate $m_N$ and $\gamma$ and then re-estimate $\beta$ using (3.95), repeating until convergence.

## 3.6 Limitations of Fixed Basis functions

Models comprising a linear combination of fixed nonlinear basis functions.

The assumption of linearity in the parameters led to a range of useful properties including closed-form solutions to the least squares problem. We can model arbitrary nonlinearities in the mapping from inputs to targets.

But there are some significant shortcomings

The basis functions $\phi_j(x)$ are fixed before the training data is observed.

The number of basis functions needs to grow rapidly with the dimensionality $D$.