Chapter 3 Linear Models for Regression Regression (supervised learning) X: D-timensional input vector t: continuous target variable Linear regression model: linear functions of adjustable parameters Linear combinations at a fixed set of nonlinear functions of input variables known as basis function. E.g.  $Y(x, w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$ 

whose basis is  $\{1, x, x^2, x^3\}$ 



Chapter 3 Linear Models for Regression

3.1 Linear Basis Function Models Linear combinations of fixed nonlinear functions of input x $y(x, w_1) := w_0 + \sum_{j=1}^{M-1} w_j \not (x)$ where  $\not (x, w)$  are known as basis functions. We allows for

any fixed offset, called bias. So it is convenient to define an additional dummy

basis function  $\varphi_o(x) = 1$  so that

$$Y(\mathcal{K}, W) = \sum_{i=1}^{M-1} W_i \mathcal{O}_i(\mathcal{K}) = W^{\mathsf{T}} \Phi(\mathcal{K})$$
  
Chapter 3 Linear Models for Begression

where  $W = (W_0, ..., W_{M-1})^T$  and  $\overline{\Phi} = (\underline{\phi}_0, ..., \underline{\phi}_{M-1})^T$ In view of pre-processing or feature extraction, the feature can be expressed as {\$\$; C\$\$)} Basis functions In Chapter 1. there is a single input x and the basis  $x ( \varphi x) = x^{i}$ functions take the form of powers of One limitation of polynomial basis: global functions

## Gaussian basis functions

$$\phi_{j}(x) = \exp\{-\frac{(x-M_{j})^{2}}{2s^{2}}\}$$

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

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

Equivalently we can use the 'tanh' function. Since tanh (a) = 20(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 Graussian noise model. As before, we assume that target t is given by a deterministic function y(x,w) with additive Gaussian noise  $t = \gamma(x, w) + \varepsilon$ 

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

I.e.

$$P(t|x, w, \beta) = N(t|Y(x, w), \beta^{-1}) \quad (3.8)$$

In section 1.5.5, we showed that

$$\mathbb{E}_{t} \mathbb{E}_{t} \mathbb$$

is the optimal prediction

Consider inputs 
$$X = \{X_{1}, ..., X_{N}\}$$
 with corresponding target  
 $t_{1}, ..., t_{N}$ . Let  $# := (t_{1}..., t_{N})^{T}$ . Assume  $X$  and  $#$  are  
drawn independently from (3.8). Then the likelihood function  
of  $W$  and  $\beta$  is in the form  
 $p(\#I, X, w, \beta) = \prod_{n=1}^{N} N C t_{n} I W^{T} \overline{p}(X_{n}), \beta$  (3.10)

We will drop the explicit X from expressions.  $\ln p(\# | W|, \beta) = \sum_{n=1}^{N} \ln N(\ln | W|^T \Phi(x_n), \beta^{-1})$  (3.11)

9

where SSE is defined by  $E_{p}(w) := \frac{1}{2} \sum_{n=1}^{N} \left\{ t_{n} - w^{T} \overline{\varrho} (x_{n}) \right\}^{2} \qquad (3.12)$ 

Consider first the maximization of (3.11) writ W. Maximization of likelihood function under a conditional Graussian (3.10) for a linear model (=> Minimizing ED CW)  $\nabla_{w_1} \ln p(\#|w_1,\beta) = \beta \sum_{n=1}^{N} (t_n - w^T \Phi(x_n)) \Phi(x_n)$ 

**Chapter 3 Linear Models for Regression** 

Here  $\Phi$  is an NXM matrix called design matrix

are given by  $\overline{\Phi}_{n;} = \varphi_j (x_n)$ 

elements

whose

I.e.

$$\overline{\Phi} = \begin{pmatrix} \varphi_{0}^{c}(\mathscr{K}_{1}) & \varphi_{1}^{c}(\mathscr{K}_{1}) & \cdots & \varphi_{M-1}^{c}(\mathscr{K}_{1}) \\ \varphi_{0}^{c}(\mathscr{K}_{2}) & \varphi_{1}^{c}(\mathscr{K}_{2}) & \varphi_{M-1}^{c}(\mathscr{K}_{2}) \\ \vdots & \vdots \\ \varphi_{0}^{c}(\mathscr{K}_{N}) & \vdots & \varphi_{M-1}^{c}(\mathscr{K}_{N}) \end{pmatrix} = \begin{pmatrix} \overline{\Phi}^{c}(\mathscr{K}_{1})^{T} \\ \overline{\Phi}^{c}(\mathscr{K}_{2})^{T} \\ \vdots \\ \overline{\Phi}^{c}(\mathscr{K}_{N})^{T} \end{pmatrix}$$

where 
$$\Phi(x) = (\phi_{0}(x), \dots, \phi_{M-1}(x))^{T}$$
.

The quantity

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

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

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

Let us see the role of bias parameter wo

$$E_{P}(w_{r}) = \frac{1}{2} \sum_{n=1}^{N} \left| t_{n} - w_{o} - \sum_{j=1}^{M-1} w_{j} \phi_{j} (x_{n}) \right|^{2}$$
 SSE

Set 
$$\frac{\partial E_{D}}{\partial w_{0}} = 0$$
 and solving for  $w_{0}$ . Then we obtain  
 $w_{0} = \overline{E} - \sum_{j=1}^{N-1} w_{j} \overline{\phi}_{j}$ 

where we have defined  $\overline{t} := \frac{1}{N} \Sigma t_n$ ,  $\overline{\phi}_j := \frac{1}{N} \Sigma \phi_j (\mathcal{X}_n)$ 

Thus, we is the difference between the averages of the and the weighted sum Chapter 3 Linear Models for Regression basis function values 14

After finding 
$$W_{ML}$$
, we can maximize log likelihood (3.11)  
w.r.t noise precision parameter  $\beta$ ,  
$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^{N} \left[ t_n - W_{ML}^T \not(X_n) \right]^2$$
residual variance of target

\_

## 3.1.2 Geometry of Least squares

N-dim

 $arphi_1$ 

 $arphi_2$ 

4 = (1)

**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 **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  $\phi_j(\mathbf{x}_n)$ .

N - dimensional space whose axes are given by tn. For fixed j, basis function values  $p_j(x_n)$  (N data points) can be represented as a vector in the same space Chapter 3 Linear Models for Regression N - dim 16 Denote  $P_{j}$  by this vector, given by  $P_{j} := (P_{j} c_{x_{1}}), P_{j} c_{x_{2}}), \dots, P_{j} c_{x_{N}})^{T}$  (  $j^{\text{th}}$  colum of  $\overline{\Phi}$ )

where 
$$j = 0, 1, ..., M-1$$
.  
Let  $M < N$  and  $5$  be the  $M$ -dim subspace spanned by  $\mathcal{P}_{j}$ .  
Define  $\mathcal{Y} := (\mathcal{Y}(\mathcal{X}_{1}, \mathcal{W}), \mathcal{Y}(\mathcal{X}_{2}, \mathcal{W}), ..., \mathcal{Y}(\mathcal{X}_{N}, \mathcal{W}))^{T}$ . Because  $\mathcal{Y}$   
 $= \sum_{j=0}^{\infty} \mathcal{W}_{j} \mathcal{P}_{j}(\mathcal{X}_{1})$   $\mathcal{Y} = \mathcal{W}_{0} \mathcal{P}_{0} + \mathcal{W}_{1} \mathcal{P}_{1} + \cdots + \mathcal{W}_{M-1} \mathcal{P}_{M-1}$   
is an arbitrary linear combination of  $\mathcal{P}_{j}$ ,  $\mathcal{Y}$  live anywhere  
in the  $M$ -dimensional subspace  $S$ .

SSE (3.12) is equal (up to a stactor  $\frac{1}{2}$ ) to  $\| y - t \|^2$ (squared Euclidean distance) Thus the least square solution w corresponds to that choice of y lying in subspace S and that is closest to t. => y = Prois(t) and  $y = \overline{\Phi} w_{ML}$ 

3.1.3 Sequential learning Aka on-line algorithm Applying the technique of stochastic gradient descent, also as sequential gradient descent Known If the error function  $E = \sum_{n} E_{n}$ , then after presentation of pattern n, SGTD updates W/ using  $W' := W' - 7 \nabla E_n$ where C: iteration number, 7: learning rate parameter.

For the case of SSE (3.12),

$$W^{(CC+1)} = W^{(C)} + \eta (t_n - W^{(C)T} \overline{f_n}) \overline{f_n}$$

where 
$$\overline{\Phi}_n := \overline{\Phi}(\mathcal{K}_n) = (\phi_0(\mathcal{K}_n), \phi_1(\mathcal{K}_n), \dots, \phi_{n-1}(\mathcal{K}_n))^T$$

$$E_{P}$$
 cm) +  $\gamma E_{W}$  cm)

where 
$$7$$
 is the regularization coefficient.  
Simple for of regularizer is as follow  
 $E_{W}(W) := \frac{1}{2} W^{T} W$  weight decay  
parameter shrint age

If we also consider SSE, then the total error function

becomes 
$$\frac{1}{2} \sum_{n=1}^{N} (t_n - w^T \Phi (x_n))^2 + \frac{7}{2} w^T w$$
 (3.27)

Set the gradient of (3.20) w.r.t W/ to zero and solve for Wr. Then we obtain the solution W/

 $W' = (\pi I + \overline{\Phi}^T \overline{\Phi})^T \overline{\Phi} \#$ 

More general regularizer is used as follows

$$\frac{1}{2} \sum_{n=1}^{N} \left\{ t_n - w^T \Phi(x_n) \right\}^2 + \frac{7}{2} \sum_{j=1}^{M} |w_j|^q \qquad (3.2q)$$

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) (=) Minimize  $E_D(W)$  subject to the constraints  $\sum_{j=1}^{M} |w_j|^q \leq \gamma$ 

Sor some appropriate Chapter 3 Linear Models for Regression

$$X \in \mathbb{R}^{P}$$
 in put, basis function  $\Phi(\cdot) = \begin{pmatrix} 1 \\ \#(c^{*}) \\ \vdots \\ \#_{M-1}(*) \end{pmatrix}$ 

- -

target 
$$t$$
, determine  $w$   
 $1 - dim$  Graussian  
 $\gamma(x, w) = w^{T} \overline{\varphi}(x) \sim t \quad (t - N(t) \gamma(x, w), \beta)$ 

target 
$$\# (K - \dim)$$
, determine  $W$   
 $(K - \dim) = W^T \Phi(K) - \# (H - N(H) \Psi(K, W), \beta I)$ 

Suppose conditional distribution of the target vector to be an istropic Gaussian  $p(\#1 \times , W, (3) = N C \# 1 W/^T \notin C \times), (3^T I)$ 

Given  $k - \dim N$  observations  $t_1, t_2, ..., t_N$ , we can combine these into  $N \times k$  matrix T. Similarly combine the input vectors  $x_1, x_2, ..., x_N$  into  $N \times D$ matrix X. The log likelihood

$$p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) = \sum_{n=1}^{N} p(\mathbf{t}_n | \mathbf{W}^{\mathsf{T}} \Phi(\mathbf{x}_n), \beta^{\mathsf{T}} \mathbf{I})$$

$$= \frac{N k}{2} l_n \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^{N} || t_n - W|^T \bar{\varrho} (m) ||^2$$

Maximization solution for W/ is given by  $M_{MX} = (\overline{\Phi}^T \overline{\Phi})^{-1} \overline{\Phi}^T T \qquad M \times k \text{ matrix}$   $M_{ML} = (\overline{\Phi}^T \overline{\Phi})^{-1} \overline{\Phi}^T T \qquad M \times k \text{ matrix}$ 

If we examine this result for each target variable 
$$t_k$$
,  
basis function values of observations  
 $W_k = (\overline{\Phi}^T \overline{\Phi})^T \overline{\Phi}^T \#_k = \overline{\Phi}^T \#_k$   
 $W_k = (\overline{\Phi}^T \overline{\Phi})^T \overline{\Phi}^T \#_k = \overline{\Phi}^T \#_k$ 

where  $\#_k$  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.

3.2 The Bias - Variance Decomposition  
Frequentist view at model complexity  
Bias - Variance trade-off  
When error function is SSE, the optimal prediction is given  

$$by$$
  $h(x) = E[t(x] = \int t p(t(x)) dt$ 

We showed in Section 1.5.5 that the expected squared loss  
can be written in the form  
Prediction optimal solution
$$(3.37)$$

$$E[L] = \int \frac{1}{7}(x_{x}) - h(x_{x}) \int_{-1}^{2} p(x_{x}) dx_{x} + \int \frac{1}{7} h(x_{x}) - t^{2} \frac{1}{7} p(x_{x}, t) dx_{x} dt$$
Chapter 3 Linear Models for Regression

$$E[L] = \int \frac{1}{2} \frac{1$$

The second term arises from the intrinsic noise and the minimum expected loss. is The first term depends on our choice of YCX) is to seek Y(x) making the first term a minimum. Our goal Culoterel & P q & il ggenze or. But goese N mel de if D  $W \rightarrow Y(x, w)$ D 3 modeling hox) PN 2N VES using y(x, w)

Bayesian: uncertainty is expressed through a posterior distribution over W/ estimate of W point based on D. Frequentist : observations D are independently drawn  $\mathcal{N}$ from P (+\*) For a given D, we can obtain a prediction function YC\*;D) Y(\*iD) and its squared error depend on D of learning algorithm is assessed The performance by taking the average over ensemble of data sets

Consider the first term in (3.37)

which depends on D.

$$f y(x; D) \pm E_D [Y(x; D)] - h(x) f^2$$

$$= \left\{ y(x; p) - \mathbb{E}_{p} [y(x; p)] \right\}^{2} + \left\{ \mathbb{E}_{p} [y(x; p)] - h(x) \right\}^{2}$$

+ 
$$2 \{ Y(X; D) - E_{P} [Y(X; D)] \{ E_{D} [Y(X; D)] - h(X) \}$$

Take the expectation w.r.t D.  

$$E_{p}[\{y(x; D) - h(x)\}^{2}] = \{E_{p}[y(x; D)] - h(x)\}^{2}$$

$$= \{E_{p}[\{y(x; D) - E_{p}[y(x; D)]\}^{2}\}$$

$$= E_{p}[\{y(x; D) - E_{p}[y(x; D)]\}^{2}]$$

$$= Variance$$

bias: extent to which the average prediction over all data sets differs from the desired regression function variance: extent to which the solutions for individual data sets vary around thier average. (sensitivity Stimper Models for Regression the choice of D)<sup>33</sup>

expected squared loss = (bias)<sup>2</sup> + variance + noise

where 
$$(bias)^2 = \int \{E_D [Y(x; D)] - h(x)\}^2 P(x) dx$$

variance = 
$$\int \mathbb{E}_{D} \left[ \frac{1}{Y} (x; D) - \mathbb{E}_{D} \left[ \frac{1}{Y} (x; D) \right] \right] p(x) dx$$
  
noise =  $\int \frac{1}{h} (x) - t \int_{x}^{2} p(x, t) dx dt$ 



Examine the bias - variance trade-off quantitately L prediction models  $y^{(2)}$ , l=1, ..., LThe average prediction  $D_1 \dots D_L$ 

$$\overline{y}(x) := -\frac{1}{L} \sum_{e=1}^{L} y^{(e)}(x)$$

and integrated squared bias and integrated variance (approximated (by sum of  $x_n$ )  $(ariance = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{L} \sum_{n=1}^{L} \{y^{(L)}(x_n) - \overline{y}(x_n)\}^2$


So the corresponding conjugate prior is given by pcwr) := N cwr ( mb, So)

.

where mean 
$$m_{o}$$
, covariance  $S_{o}$ .  
Thus the posterior distribution in the form (see (2.116))  
 $p(W|H) = N(W|M_{N}, S_{N})$   
where  $m_{N} = S_{N}(S_{o}^{-1}m_{b} + \beta \Phi^{T} \pm)$ 

 $S_{N}^{-1} = S_{0}^{-1} + \beta \overline{\Phi}^{T} \overline{\Phi}$ 

Since the posterior is Graussian (unimodal), its mode = mean  $\Rightarrow W_{MAP} = MI_N$ 

If we consider  $S_0 := \alpha^2 I$  and  $\alpha \rightarrow 0$  i.e. infinitely broad prior, then the mean  $M_N$  reduces to  $M_{ML}$ 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  
 $P(W|\alpha) = N(W|0, \alpha^{-1}I)$  (simple version)

So the corresponding posterior

$$pcw(#) = N cw/ m_N, S_N)$$

where

$$m_{N} = \beta S_{N} \overline{\Phi}^{T} \overline{\Phi}$$
$$S_{N}^{-1} = \alpha I + \beta \overline{\Phi}^{T} \overline{\Phi}$$

Log of posterior is the sum of log of litelihood and log of prior  $\ln p(w(1t)) = -\frac{\beta}{2} \sum_{n=1}^{N} (t_n - w)^T \overline{p}(w_n))_1^2 - \frac{\alpha}{2} w_n^T w_n + constant$ 

Its MAP solution w.r.t W' is equivalent to minimization of SSE with additional quadratic regularization term  $7 = K/\beta$ .



- input: x
- target: t
- $y(x, w) := w_0 + w_1 x$
- Observed data o generated by -0.3 + 0.5 x with std 0.2



Figure for the second 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 Graussian prior  

$$p(wr(\alpha)) := \left[\frac{q}{2}\left(\frac{\alpha}{2}\right)^{kq} \frac{1}{P(l'q)}\right]^{M} \exp\left(-\frac{q}{2}\sum_{j=0}^{M-1}|w_{ij}|^{q}\right)$$

in which q=2 corresponds to the Graussian.

If q=2, MAP solution of wr is the minimization solution

of (3.29) which is SSE t regularization term

If q ≠ 2 it is not true. (mode of posterior ≠ mean)

3.3.2 Predictive distribution  
In practice, we are interested in making predictions of t  
for a new \* ( not the value of W)  
Predictive distribution of t  

$$p(t|*, #, \alpha, \beta) = \int p(t|*, w, \beta) p(w|t, \alpha, \beta) dw$$
  
new input  
where # is the vector of training target values.  
 $\alpha$  is from prior assumption,  $\beta$  is Gaussian noise of t  
 $p(t|*, w, \beta) = N(t|Y(*, w), \beta^{-1})$   
Chapter 3 Linear Models for Regression

The predictive distribution takes the sorm  
, training data 
$$\in \mathbb{R}$$
  
 $P(t| \times, \#, \propto, \beta) = \mathcal{N}(t| m_{\mathcal{N}}^{T} \Phi(x), \sigma_{\mathcal{N}}^{2}(x))$   
(new input

where

$$\vec{\sigma}_{\mathcal{N}}(\mathcal{K}) = \frac{1}{\beta} + \vec{\rho} (\mathcal{K})^{\mathsf{T}} S_{\mathcal{N}} \vec{\phi} (\mathcal{K}) \quad (3.59)$$

$$\text{data noise} \quad \text{uncertainty}$$

$$\text{of } \mathcal{M}$$

By [Oazaz et al., 199n] $\sigma_{NH}(x) \leq \sigma_{N}^{2}(x)$ 

If N + 00, then Chapter 3 Linear Models for Regression (3.59) - 1/3 45



**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 **Bapter**ical interview **Constant Support Constant <b>Constant Constant Constant Con** 



**Figure 3.9** Plots of the function y(x, w) using samples from the posterior distributions over 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 β<sup>-1</sup>.
  - W posterior a sampling  $W_{1...} W_{M}$  $\int \int dw = \frac{1}{M} \sum (w_{n})$

3.3.3 Equivalent kernel (kernel method)  
Substitute (3.53) into (3.3) (expected prediction)  

$$\begin{array}{l} \text{MX1} \\ \text{Y(X, M_N)} = \text{M}_N^T \overline{\Phi}(X) = \beta \overline{\Phi}(X)^T S_N \overline{\Phi}^T = \sum_{n=1}^N \beta \overline{\Phi}(X)^T S_N \overline{\Phi}(X) + n \\ \text{MX1} \end{array}$$
where  $\overline{\Phi}(X) := (\beta (X) \dots \beta_{MT}(X))^T$ ,  $S_N^{-1} = S_0^{-1} + \beta \overline{\Phi}^T \overline{\Phi}$  and  $\overline{\Phi} = \begin{pmatrix} g(X) \dots g_{MT}(X) \\ \vdots \dots \\ g(X) \end{pmatrix} \text{ design matrix}$ 

$$N \times M$$

Thus, YCX, MN, is the linear combination of the training. set target variables <sup>Chapter 3</sup> Linear Models for Regression 49

$$Y(\mathcal{H}, \mathcal{M}_{\mathcal{N}}) = \sum_{n=1}^{\mathcal{N}} k(\mathcal{H}, \mathcal{H}_n) t_n$$

where the Sunction

 $\Rightarrow$ 

$$F(*,*') := \beta \Phi(*)^T S_v \Phi(*')$$

is known as smoother matrix or <u>equivalent bernel</u> Linear smoother: regression function makes predictions by taking linear combinations of training target values This kernel tepends on \*n because of S<sub>N</sub> Chapter 3 Linear Models for Regression 50 Consider the covariance between YCXY) and YCXY)

$$cov[y(x), y(x')] = cov[w]^T \Phi(x), w]^T \Phi(x')]$$
$$= \Phi(x)^T S_n \Phi(x') = \beta^T k(x, x')$$

$$Y(x) = N(t|m_{\mu}^{T} \overline{\mathcal{Q}}(x), \beta^{T} + \overline{\mathcal{Q}}(x)^{T} S_{\mu} \overline{\mathcal{Q}}(x))$$
 (not scalar value)

 $y(x) = w^T \Phi(x)$  where  $W = N(W|M_N, S_N)$ 

 $cov \ \Box \ w^{T} \ \overline{\Phi} (x), \ w^{T} \ \overline{\Phi} (x') \ ] = \underbrace{\mathbb{E} \left[ \ \overline{\Phi} (x)^{T} \ w \ w^{T} \ \overline{\Phi} (x') \right] - \Phi (x)^{T} \ m_{w} \ m_{w}^{T} \ \overline{\Phi} (x')}_{= \overline{\Phi} (x)^{T} \ \overline{\mathbb{E}} \left[ w \ w^{T} \ \overline{\Phi} (x') \right] + \overline{\mathbb{E}} \left[ w \ w^{T} \ \overline{\Phi} (x') \right]^{T}}_{= cov \ [w] \ + \ \mathbb{E} \left[ w \ w^{T} \ \overline{\Phi} (x') \right]^{T}}$ 

Chapter 3 Linear Models for Regression Min Min

51

For	regres	sion,	we int	roduces	d a	set	र्व	basis	funct	ions	SO
equiv	alent	kerne	l was	impl	icitly	<b>4e</b> t	ermina	ઝ.			
But	we	con	define	0	locali	हरू	kem	e (	firectly	and	use
this	to	make	predict	znoś							

The equivalent kernel (3.62) can be expressed in the form on inner product w.r.t  $\Psi(x)$  of nonlinear functions  $\kappa(x, \overline{x}) := \Psi(x)^T \Psi(\overline{x})$ 

where 
$$\underline{\Psi}(x) := \beta^{1/2} S_{N}^{1/2} \overline{\Phi}(x)$$

3.4	Bayesian Model		Model	Compariso						
The	problem	ঠ	model	selection	from	a	Bayesian	pers	pecti	ve.
The	over - fit	<i>iting</i>	مجمحذم	ted with	Maxi	mam	likeliho	od ca	n b	e
avoided by marginalizing over model parameters.										
The	Bayesian		to us	model	compan	ison	involves	the	use	of
proba	bilities	to	represe	nt unc	ertointy	ir	the c	choice	र्व	model.

Compare a set of L models { Milizin.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).



**Chapter 3 Linear Models for Regression** 

This uncertainty is expressed through  $P(M_{i})$ a training data set D we want to evaluate Given 9 PCM; ID) oc PCM; ) P(DIM; ) PCD(O)prior model evidence The prior can express a preference for different models. But for simplicity assume that all models have the same prior PCPIMi) 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 morginalized out Models for Regression

56

The predictive distribution (mixture distribution) posterior over models  $P(t | X, D) = \sum_{i=1}^{L} P(t | X, M_i, D) P(M_i, D)$ 

Average of the predictive distributions P(t | \*, Mi, D) of individual models weighted by the posterior probabilities p(M; ID) For example, two models M2 M, 0 Model selection: use the single most probable model alone.

Consider model Mi governed by the parameter w. The model evidence is given by  $P(P|M_i) = \int P(P|W, M_i) P(W|M_i) dW$ 

The model evidence (marginal likelihood) pCDIM;) 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

$$P(W|D, M_i) = \frac{P(D|W, M_i) P(W|M_i)}{P(D|M_i)}$$

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



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  $ln p(D) \simeq ln p(D|w_{MAP}) + ln \left(\frac{\Delta W_{posterior}}{\Delta W_{prior}}\right)$   
 $fit to dote given by the most probable complexity
 $\Delta W_{posterior} < \Delta W_{prior}$ , 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 tota in  
posterior, then the Charter Builders of Regression arge.$ 

.

For a model with M parameters, assume all parameters have  
the same ratio of 
$$\Delta W_{\text{posterior}} / \Delta W_{\text{prior}}$$
, then we obtain  
a similar approximation as follows  
 $\ln P(D) \simeq \ln P(D | W_{\text{MAP}}) + M \ln \left(\frac{\Delta W_{\text{posterior}}}{\Delta W_{\text{prior}}}\right)$ 

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 & and B - Make predictions by marginalizing w.r.t these hyperparameters and parameters W/ the complete marginalization over all of these variables But X.B and w/ is analytically intractable.



**Chapter 3 Linear Models for Regression** 

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  $posterior$  over w  
 $p(t 1 \pm) = \int \int p(t + w, \beta) p(w + \pm, \alpha, \beta) p(\alpha, \beta + \pm) dw d\alpha d\beta$   
model assumption  
(3.8)  
Here we omitted the dependence on input  $\%$ .

If posterior pcx, ß (++) is sharply peaked around à and ß, then

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

Bayes Theorem, the posterior distribution for X, B From PCX, B(#) oc P(# | X, B) P(X, B) is relatively flat, the values & and B are So if prior obtained by maximizing the marginal likelihood function plt (x,p) **Chapter 3 Linear Models for Regression** 

Here, we evaluate the marginal litelihood for the linear basis then finding its model maxima. and us to determine values for hyperparameters So this will allow training tota alone ( «/3 = regularization parameter) from the approaches of maximization of the log evidence Tuo the evidence function analytically and then - Evaluate set its derivative O to obtain re-estimation to equal for X, B - Use the technique called expectation maximization algorithm in Section 9.3.4.

3.5.1 Evaluation of the evidence function The marginal likelihood function pctla, B) is obtained by integrating over w

$$P(\# | X, \beta) = \int P(\# | W, \beta) P(W | X) dW$$

By the result (2.115) for the conditional distribution in a linear - Graussian 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)

$$P(\# \mid \alpha, \beta) = \left(\frac{\beta}{2\pi}\right) \left(\frac{\alpha}{2\pi}\right) \int \exp\{-E(w)\} dw \qquad (3.78)$$

where M is the dimensionality of W and  

$$E(W) = \beta E_D(W) + \alpha E_W(W)$$

$$= \frac{\beta}{2} \| \# - \overline{\Phi} W \|^2 + \frac{\alpha}{2} W^T W$$
(3.79)

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

where we have introduced

$$A = \kappa I + \beta \overline{\Phi}^T \overline{\Phi}$$

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

 $A = \nabla \nabla E C W$ 

$$M_{N} = \beta A^{-1} \overline{\Phi}^{T} \overline{E}. \qquad (3.84)$$

 $S_{N}^{-1} = KI + \beta \Phi^{T} \Phi$ Using (3.54), we see  $A = S_{N}^{-1}$ , hence (3.84) = (3.53)

Back to the integration (3.18)  

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

$$= \exp\{-E(m_W)\} \int \exp\{-\frac{1}{2}(w - m_W)^T A(w - m_W)\} dw$$

$$= \exp\{-E(m_W)\} (2\pi)^{M_2} |A|^{-N_2}$$
Using (3.18), we can write the log of the marginal  
likelihood in the form  

$$\ln p(\#(x,\beta) = \frac{M}{2} \ln x + \frac{N}{2} \ln \beta - E(m_W) - \frac{1}{2} \ln |A| - \frac{N}{2} \ln (2\pi)$$



Remark

- The underlying sinusoidal function is an odd function - In M=3 case, we obtain a significant improvement in

data Sit.
3.5.2 Maximizing the evidence function Consider the maximization of P(t(x, b) w.r.t x. This can be done by first defining the following eigenvector equation

$$(\rho \Phi^{T} \Phi) u_{\lambda} = \eta_{\lambda} u_{\lambda}$$

Since  $A = \alpha I + \beta \overline{\Phi}^T \overline{\Phi}$ , A has eigenvalues  $\alpha + \pi_i$ Now consider the partial derivative of ln |A| w.r.t  $\alpha$ .

**Chapter 3 Linear Models for Regression** 

$$\frac{\partial}{\partial \alpha} \ln[A] = \frac{\partial}{\partial \alpha} \ln \prod_{i} (n_{i} + \alpha) = \frac{\partial}{\partial \alpha} \sum_{i} \ln(n_{i} + \alpha) = \sum_{i} \frac{1}{n_{i} + \alpha}$$

Thus

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

Multiplying by 2x and rearranging, we obtain  $\chi m_N^T m_N = M - \chi \sum_{\lambda} \frac{1}{7_{\lambda} + \chi} =: r$  Since there are M terms in the sum over  $\lambda$ , (3.91)  $f = \sum_{\lambda} \frac{\eta_{\lambda}}{\alpha + \eta_{\lambda}}$  (depends on  $\alpha$ )

So the following x maximizes the marginal litelihoot  $x = \frac{t}{m_N^T m_N}$  (3.92)

Note that i depends on a and the mode  $M_N$  of the posterior distribution depends on the choice of a. Thus, this solution is implicit and is adopted an iterative procedure Chapter 3 Linear Models for Regression 75

Make an initial choice of x and use this to find  $m_{l_{A}}(3.53)$  and evaluate f(3.91)Using (3.92), re-estimate & and the process repeat until convergence. Note that because the matrix  $\overline{\mathcal{D}}^T \overline{\mathcal{D}}$  is fixed, we can compute its eigenvalues once at the start The value of & has been determined purely by training tata.

Similarly, maximize the log marginal likelihood (3.86) w.r.t (3)  
Note that the eigenvalues 
$$7\lambda$$
 are proportional to  $\beta$   
and hence  $\frac{d^{2}\lambda^{2}}{d\beta} = \frac{7\lambda^{2}}{\beta}\beta$  giving  
 $\frac{\partial}{\partial\beta}\ln|A| = \frac{\partial}{\partial\beta}\sum_{\lambda}\ln(7\lambda + \alpha) = \frac{1}{\beta}\sum_{\lambda}\frac{7\lambda}{7\lambda + \alpha} = \frac{d}{\beta}$ 

So the stationary point of the marginal likelihood satisfies  $0 = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^{N} \beta t_n - m_N^T \Phi(x_n) \beta^2 - \frac{1}{2\beta}$ 

and rearranging we obtain

$$\frac{1}{\beta} = \frac{1}{N-\delta} \sum_{n=1}^{N} \left\{ t_n - m \right\}_{N}^{T} \overline{\Phi} \left( x_n \right) \right\}^{2}$$

Again, this is an implicit solution for  $\beta$ . So choose an initial value for  $\beta$  and calculate  $M_N$  and it 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 least squares problem. We can model arbitrary nonlinearities the the mapping from inputs to targets. in But there are some significant short comings

The basis functions \$(1) are fixed before the training data is observed.

The number of basis functions needs to grow rapidly

with the dimensionality D.