

PRML Note

C03 Linear Models for Regression

Yang Zhao

Department of Automation, Tsinghua University

- Given a training data set comprising N observations $\{\mathbf{x}_n\}$, together with corresponding target values $\{t_n\}$, the goal is to predict the value of t for a new value of \mathbf{x} .
- We can not only use the linear functions of the input variables, but also use the linear combinations of a fixed set of nonlinear functions of the input variables, known as *basis functions*. Such models are linear functions of the parameters and can be nonlinear with respect to the input variables.

1 Linear Basis Function Models

- The simplest linear model for regression is

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \cdots + w_Dx_D$$

We can extend it by considering linear combinations of fixed nonlinear functions of the input variables, of the form

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

where $\phi_j(\mathbf{x})$ are known as basis functions. We can also define an additional dummy basis function $\phi_0(\mathbf{x}) = 1$ so that

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (2)$$

where $\mathbf{w} = (w_0, \cdots, w_{M-1})^T$ and $\boldsymbol{\phi} = (\phi_0, \cdots, \phi_{M-1})^T$.

- *polynomial basis function*

$$\phi_j(x) = x^j$$

One limitations of this basis functions is that they are global functions of the input variables, so that changes in one region of input space affect all other regions. This can be resolved by dividing the input space up into regions and fitting a different polynomial in each region, leading to *spline functions*.

- *Gaussian basis function*

$$\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}} \quad (3)$$

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

- *sigmoidal basis function*

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right) \quad (4)$$

where σ is the logistic sigmoid function defined by

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad (5)$$

- We assume that the target variable t is given by a deterministic function $y(\mathbf{x}, \mathbf{w})$ with additive Gaussian noise so that

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

where ϵ is a zero mean Gaussian random variable with precision β . Thus

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (6)$$

- Consider a data set of inputs $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ with corresponding target values $\mathbf{t} = (t_1, t_2, \dots, t_N)^T$, which are drawn independently from the distribution 6, so we have

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

The gradient of the loglikelihood function takes the form

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) \quad (7)$$

Setting this gradient to zero gives

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n) - \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right) \mathbf{w}$$

Solving for \mathbf{w} , we obtain

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (8)$$

where

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

- Considering the w_0 , rewritten the error function, we have

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n))$$

Setting the derivative with respect to w_0 equal to zero, and solving for w_0 , we obtain

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

where we define

$$\bar{t} = \frac{1}{N} \sum_{n=1}^N t_n$$

$$\bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$$

Thus the bias w_0 compensates for the difference between the averages(over the training set) of the target values and the weighted sum of the averages of the basis function values.

- We can also maximize the log likelihood function with respect to the noise precision parameters β , giving

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n))^2$$

- When $\boldsymbol{\Phi}^T \boldsymbol{\Phi}$ is close to singular, solving this problem directly will have numerical difficulties. We can use *singular value decomposition*, or SVD to address this difficulties.
- *LMS algorithm*. This is also known as *least-mean-squares*. By applying the technique of *stochastic gradient descent*, we can update the \mathbf{w}

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \nabla E_n \\ &= \mathbf{w}^{(t)} - \eta (\mathbf{w}^{(t)T} \boldsymbol{\phi}_n - t_n) \boldsymbol{\phi}_n \end{aligned}$$

we can see how it works in figure 1

- By adding a regularization term to error function in order to control over-fitting, we get

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

This choice of regularizer is known in the ML literature as *weight decay*. We obtain

$$\mathbf{w} = (\lambda \mathbf{I} + \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \quad (9)$$

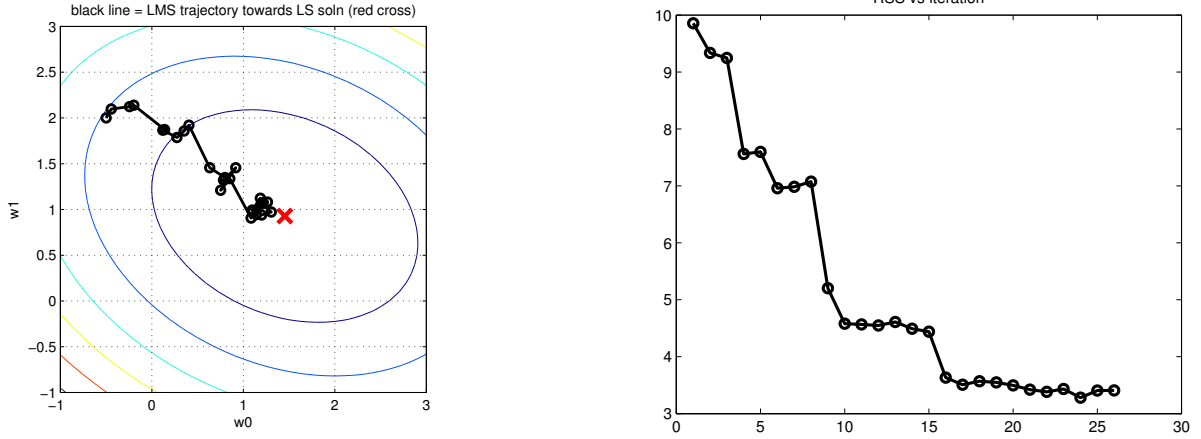


Figure 1: Illustration of LMS algorithm. Note that it doesn't decrease monotonically

- The problem of determining the optimal model complexity is then shifted from finding the appropriate number of basis functions to determining a suitable value of the regularization coefficient λ .
- For multiple outputs, we could solve this problem by introducing a different set of basis functions for each component of \mathbf{t} . But a more interesting and more common approach is to use the same set of basis functions to model all of the components of target vector, so that

$$y(\mathbf{x}, \mathbf{W}) = \mathbf{W}^T \phi(\mathbf{x}) \quad (10)$$

- Consider a linear basis function regression model for a multivariate target variable \mathbf{t} having a Gaussian distribution of the form

$$p(\mathbf{t}|\mathbf{W}, \Sigma) = \mathcal{N}(\mathbf{t}|\mathbf{y}(\mathbf{x}, \mathbf{W}), \Sigma) \quad (11)$$

where $\mathbf{y}(\mathbf{x}, \mathbf{W}) = \mathbf{W}^T \phi(\mathbf{x})$, so the log-likelihood function is

$$\ln p(\mathbf{T}|\mathbf{W}, \Sigma) = \sum_{i=1}^N (\mathbf{t}_i - \mathbf{W}^T \phi(\mathbf{x}_i))^T \Sigma^{-1} (\mathbf{t}_i - \mathbf{W}^T \phi(\mathbf{x}_i))$$

where

$$\mathbf{T} = \begin{pmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \\ \vdots \\ \mathbf{t}_N \end{pmatrix} = \begin{pmatrix} t_{11} & t_{12} & \cdots & t_{1K} \\ t_{21} & t_{22} & \cdots & t_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ t_{N1} & t_{N2} & \cdots & t_{NK} \end{pmatrix}$$

So, let

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} \ln p(\mathbf{T}|\mathbf{W}, \Sigma) &= \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \mathbf{W} \Sigma^{-1} - \phi(\mathbf{x}_i) \mathbf{t}_i \Sigma^{-1} \\ &= \Phi^T \Phi \mathbf{W} \Sigma^{-1} - \Phi^T \mathbf{T} \Sigma^{-1} \\ &= \mathbf{0} \end{aligned} \quad (12)$$

we can get

$$\mathbf{W} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$$

2 The Bias-Variance Decomposition

- Although the introduction of regularization terms can control over-fitting for models with many parameters, this raises the question of how to determine a suitable value for the regularization coefficient λ .
- Using an unbiased estimator is not always the best according to the bias-variance trade-off.
- For the squared loss function, the optimal prediction is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x})dt$$

Accordingt to Chapter 1, we known that

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \iint (h(\mathbf{x}) - t)^2 p(t, \mathbf{x}) dt d\mathbf{x} \quad (13)$$

The second term is independent of $y(\mathbf{x})$, representing the minimum achievable value of the expected loss. For any given data set \mathcal{D} , we can get a prediction function $y(\mathbf{x}; \mathcal{D})$, and the average over the different data sets denotes by $\mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}))$. So the first term in the equation 13 can be rewritten by

$$\begin{aligned} \text{the first term} &= \int (y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} \\ &= \int (y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D})) + \mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}))^2) p(\mathbf{x}) d\mathbf{x} \\ &= \{\mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D})) - h(\mathbf{x})\}^2 + \mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}\{y(\mathbf{x}; \mathcal{D})\})^2 \end{aligned} \quad (14)$$

- the first term in equation 14 is called *the squared bias*, which represents the extent to which the average prediction over all data sets differs from the desired regression.
- the second term, called the *variance*, measures the extent to which the solutions for individual data sets vary around their average.

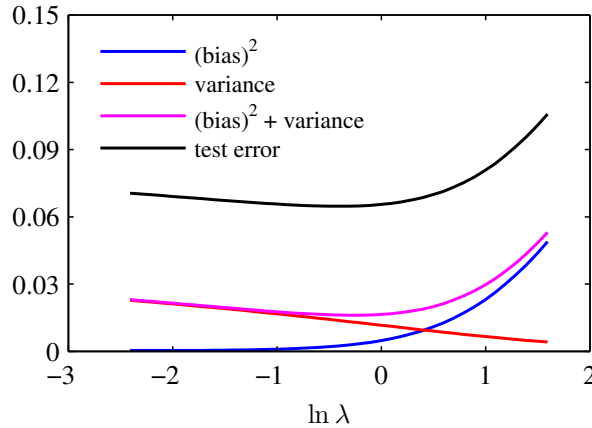


Figure 2: The bias-variance decomposition

- The model with the optimal predictive capability is the one that leads to the best balance between bias and variance.
- The bias-variance decomposition is based on averages with respect to ensembles of data sets, whereas in practice we have only the single observed data set. And the larger data set can support the larger capability of the model. So splitting the data set into many sets is not good for most situations.

3 Bayesian Linear Regression

- Given the conjugate prior

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$$

and the likelihood function

$$p(\mathbf{t}|\mathbf{w}) = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

So we can write the posterior distribution directly in the form

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N)$$

where

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta\Phi^T\Phi\end{aligned}$$

- If data points arrive sequentially, then the posterior distribution at any stage acts as the prior distribution for the subsequent data point. By completing the square in the exponential, we can obtain

$$\begin{aligned}\mathbf{S}_{N+1}^{-1} &= \mathbf{S}_N^{-1} + \beta\phi_{N+1}\phi_{N+1}^T \\ \mathbf{m}_{N+1} &= \mathbf{S}_{N+1}(\mathbf{S}_N^{-1}\mathbf{m}_N + \beta\phi_{N+1}t_{N+1})\end{aligned}$$

- Sometimes, we are not interested in the value of \mathbf{x} itself but rather in making predictions of t for new values of \mathbf{x} .
- The *predictive distribution*

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

where

$$\begin{aligned}p(t|\mathbf{x}, \mathbf{w}, \beta) &= \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \\ p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)\end{aligned}$$

So the predictive distribution takes the form

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|\mathbf{m}_N^T(\phi(\mathbf{x})), \sigma_N^2(\mathbf{x}))$$

where

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

The first term represents the noise on the data whereas the second term reflects the uncertainty associated with the parameters \mathbf{w} .

- Assume we have already observed N data points. By considering an additional $(\mathbf{x}_{N+1}, t_{N+1})$

$$\begin{aligned}\mathbf{S}_{N+1}^{-1} &= \mathbf{S}_N^{-1} + \beta \phi(\mathbf{x}_{N+1}) \phi(\mathbf{x}_{N+1})^T \\ \mathbf{m}_{N+1} &= \mathbf{S}_{N+1}(\mathbf{S}_N^{-1} \mathbf{m}_N + \phi(\mathbf{x}_{N+1}) t_{N+1})\end{aligned}$$

So by completing the square in the exponential

$$\begin{aligned}\sigma_{N+1}^2(\mathbf{x}) - \sigma_N^2(\mathbf{x}) &= \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}) - \left(\frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_{N+1} \phi(\mathbf{x}) \right) \\ &= \phi(\mathbf{x})^T (\mathbf{S}_{N+1} - \mathbf{S}_N) \phi(\mathbf{x}) \\ &= -\phi(\mathbf{x})^T \left(\frac{\beta \mathbf{S}_N \phi(\mathbf{x}_{N+1}) \phi(\mathbf{x}_{N+1})^T \mathbf{S}_N}{1 + \beta \phi(\mathbf{x}_{N+1})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1})} \right) \phi(\mathbf{x}) \\ &\leq 0\end{aligned}$$

because $\mathbf{S}_{N+1} - \mathbf{S}_N$ is negative definite matrix. So the second term in the predictive variance will go to zero.

- Considering the predictive mean, we obtain

$$y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \phi(\mathbf{x}) = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n \quad (15)$$

Thus the mean of the predictive at a point \mathbf{x} is given by a linear combination of the training set target variables t_n , so that we can write

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n \quad (16)$$

where

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n)$$

is known as the *smoother matrix* or the *equivalent kernel*, which can also be written as

$$k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})$$

where $\psi(\mathbf{x}) = \beta^{1/2} \mathbf{S}_N^{1/2} \phi(\mathbf{x})$.

- We can see that $k(\mathbf{x}, \mathbf{x}_n)$ can be weight of the linear combination and the data points close to \mathbf{x} are given higher weight than the data points further removed from \mathbf{x} .
- Note that the kernel function can be negative as well as positive, so although it satisfies a summation constraint, the corresponding predictions are not necessarily convex combinations of the training set target variables.

4 Bayesian Model Comparison

- Suppose we wish to compare a set of L models $\{\mathcal{M}_i\}$ where $i = 1, 2, \dots, L$. Given a training set \mathcal{D} , we wish to calculate the posterior distribution $p(\mathcal{M}_i)$

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i)$$

The prior allows us to express a preference for different models. Let us simply assume that all models are given equal prior probability. The interesting term is the *model evidence* $p(\mathcal{D}|\mathcal{M}_i)$ which expresses the preference shown by the data for different models. The model evidence is sometimes also called the *marginal likelihood*. The ratio of model evidence $p(\mathcal{D}|\mathcal{M}_i)/p(\mathcal{D}|\mathcal{M}_j)$ for two models is known as a *Bayes factor*.

5 The Evidence Approximation

- In a fully Bayesian treatment of the linear basis function model, we would introduce prior distributions over the hyperparameters α and β and make predictions by marginalizing with respect to these hyperparameters as well as with respect to the parameters \mathbf{w} , which is analytically intractable. We can make an approximation in which we set the hyperparameters to specific values determined by maximizing the *marginal likelihood function* obtained by first integrating over the parameters \mathbf{w} . This framework is known as *empirical Bayes*, or *type 2 maximum likelihood*, or *generalized maximum likelihood*, or *evidence approximation*.
- The predictive distribution is obtained by marginalizing over $\mathbf{w}, \alpha, \beta$, so that

$$p(t|\mathbf{x}, \mathbf{t}, \mathbf{X}) = \iiint p(t|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) p(\alpha, \beta|\mathbf{t}, \mathbf{X}) d\mathbf{w} d\alpha d\beta \quad (17)$$

where

$$\begin{aligned} p(t|\mathbf{x}, \mathbf{w}, \beta) &= \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \\ p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi \end{aligned}$$

For posterior distribution $p(\alpha, \beta|\mathbf{t}, \mathbf{X})$, we can simply assume the it is sharply peaked around values $\hat{\alpha}$ and $\hat{\beta}$, so

$$p(t|\mathbf{x}, \mathbf{t}, \mathbf{X}) \simeq p(t|\mathbf{x}, \mathbf{t}, \mathbf{X}, \hat{\alpha}, \hat{\beta}) = \int p(t|\mathbf{x}, \mathbf{w}, \hat{\beta}) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \hat{\alpha}, \hat{\beta}) d\mathbf{w}$$

and the posterior distribution for α and β is given by

$$p(\alpha, \beta|\mathbf{t}, \mathbf{X}) \propto p(\mathbf{t}|\alpha, \beta, \mathbf{X}) p(\alpha, \beta)$$

If the prior is relatively flat, the value of $\hat{\alpha}$ and $\hat{\beta}$ are obtained by maximizing the marginal likelihood function $p(\mathbf{t}|\alpha, \beta, \mathbf{X})$. We can both use the analytially method by setting the derivative equal to zero and the expectation maximization algorithm.

- The marginal likelihood function $p(\mathbf{t}|\alpha, \beta, \mathbf{X})$ is obtained by integrating over the weight parameters \mathbf{w} , so that

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

where

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{w}, \beta, \mathbf{X}) &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \\ E_D(\mathbf{w}) &= \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 \\ p(\mathbf{w}|\alpha) &= \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1} \mathbf{I}) \end{aligned}$$

- We can write the evidence function in the form

$$p(\mathbf{t}|\alpha, \beta, \mathbf{X}) = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int e^{-E(\mathbf{w})} d\mathbf{w}$$

where

$$E(\mathbf{w}) = \beta E_D(\mathbf{w}) + \alpha E_W(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

we can complete the square over \mathbf{w} giving

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

where we have introduced

$$\begin{aligned} \mathbf{A} &= \alpha \mathbf{I} + \beta \Phi^T \Phi \\ E(\mathbf{m}_N) &= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \\ \mathbf{m}_N &= \beta \mathbf{A}^{-1} \Phi^T \mathbf{t} \end{aligned}$$

So the integral over \mathbf{w} can now be evaluated as

$$\begin{aligned} \int e^{-E(\mathbf{w})} d\mathbf{w} &= e^{-E(\mathbf{m}_N)} \int e^{-\frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N)} d\mathbf{w} \\ &= e^{-E(\mathbf{m}_N)} (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \end{aligned}$$

Thus we can obtain

$$\ln p(\mathbf{t}|\alpha, \beta, \mathbf{X}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) - \frac{1}{2} \ln |\mathbf{A}| - \frac{N}{2} \ln(2\pi) \quad (18)$$

- Another way to get the equation 18. The likelihood function can be rewritten by

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

So, the result gives that

$$p(\mathbf{t}|\alpha, \beta, \mathbf{X}) = \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I}_N + \alpha^{-1} \Phi \Phi^T)$$

Taking the log we obtain

$$\ln p(\mathbf{t}|\alpha, \beta, \mathbf{X}) = -\frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln|\beta^{-1}\mathbf{I}_N + \alpha^{-1}\Phi\Phi^T| - \frac{1}{2}\mathbf{t}^T(\beta^{-1}\mathbf{I}_N + \alpha^{-1}\Phi\Phi^T)^{-1}\mathbf{t}$$

Using the equation 26, we have

$$\begin{aligned} |\beta^{-1}\mathbf{I}_N + \alpha^{-1}\Phi\Phi^T| &= \beta^{-N}|\mathbf{I}_N + \beta\alpha^{-1}\Phi\Phi^T| \\ &= \beta^{-N}|\mathbf{I}_M + \alpha^{-1}\Phi^T\Phi| \\ &= \beta^{-N}\alpha^{-M}|\alpha\mathbf{I}_M + \beta\Phi^T\Phi| \\ &= \beta^{-N}\alpha^{-M}|\mathbf{A}| \end{aligned}$$

and using the *Woodbury matrix identity* introduced in Chapter 2 appendix, we obtain

$$\begin{aligned} -\frac{1}{2}\mathbf{t}^T(\beta^{-1}\mathbf{I}_N + \alpha^{-1}\Phi\Phi^T)^{-1}\mathbf{t} &= -\frac{1}{2}\mathbf{t}^T\left(\beta\mathbf{I}_N - \beta\Phi(\alpha\mathbf{I}_M + \beta\Phi^T\Phi)^{-1}\Phi^T\beta\right)\mathbf{t} \\ &= -\frac{\beta}{2}\mathbf{t}^T\mathbf{t} + \frac{\beta^2}{2}\mathbf{t}^T\Phi\mathbf{A}^{-1}\Phi^T\mathbf{t} \\ &= -\frac{\beta}{2}\mathbf{t}^T\mathbf{t} + \frac{1}{2}\mathbf{m}_N^T\mathbf{A}\mathbf{m}_N \\ &= -\frac{\beta}{2}\|\mathbf{t} - \Phi\mathbf{m}_N\|^2 - \frac{\alpha}{2}\mathbf{m}_N^T \\ &= -E(\mathbf{m}_N) \end{aligned}$$

and we have already got the same solutions as the equation 18.

- Maximizing the evidence function 18 with respect to α . First define the following eigenvector equation

$$\beta\Phi^T\Phi\mathbf{u}_i = \lambda_i\mathbf{u}_i \quad (19)$$

Thus, we can know that the \mathbf{A} has eigenvalues $\alpha + \lambda_i$. Because

$$\begin{aligned} \frac{d}{d\alpha}\ln|\mathbf{A}| &= \frac{d}{d\alpha}\ln\prod_i(\alpha + \lambda_i) \\ &= \frac{d}{d\alpha}\sum_i\ln(\alpha + \lambda_i) \\ &= \sum_i\frac{1}{\alpha + \lambda_i} \end{aligned}$$

Thus the stationary points of equation 18 with respect to α satisfy

$$0 = \frac{M}{2\alpha} - \frac{1}{2}\mathbf{m}_N^T\mathbf{m}_N - \frac{1}{2}\sum_i\frac{1}{\lambda_i + \alpha} \quad (20)$$

So

$$\alpha\mathbf{m}_N^T\mathbf{m}_N = M - \alpha\sum_i\frac{1}{\lambda_i + \alpha} = \sum_i\frac{\lambda_i}{\lambda_i + \alpha} = \gamma \quad (21)$$

which equals to

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T\mathbf{m}_N} \quad (22)$$

Note that this is an implicit solution for α not only because γ depends on α but also because the model \mathbf{m}_N of the posterior distribution itself depends on the choice of α , which means we have to use the iterative procedure to get the optimal α .

- Maximizing the evidence function 18 with respect to β . We note that the eigenvalues λ_i is proportional to β , and hence $d\lambda_i/d\beta = \lambda_i/\beta$ giving

$$\frac{d}{d\beta} \ln|\mathbf{A}| = \frac{d}{d\beta} \sum_i (\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta} \quad (23)$$

So

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N (t_n - \mathbf{m}_N^T \phi(\mathbf{x}_n))^2 \quad (24)$$

which is also an implicit solution for β .

- The quantity γ measures the effective total number of well determined parameters.

6 Limitations of Fixed Basis Functions

- The difficulty from the assumption that the basis functions $\phi_j(\mathbf{x})$ are fixed before the training data set is observed. And the number of the basis functions needs to grow rapidly, often exponentially, with the dimensionality of the input space.
- The data vector typically lie close to a nonlinear manifold whose intrinsic dimensionality is small.
- Target variables may have significant dependence on only a small number of possible directions within the data manifold.

7 Appendix

- *The matrix identity for the sequential data*

$$(\mathbf{M} + \mathbf{v}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1}\mathbf{v})(\mathbf{v}^T\mathbf{M}^{-1})}{1 + \mathbf{v}^T\mathbf{M}^{-1}\mathbf{v}} \quad (25)$$

- If \mathbf{A} and \mathbf{B} are matrices of size $N \times M$, then

$$|\mathbf{I}_N + \mathbf{A}\mathbf{B}^T| = |\mathbf{I}_M + \mathbf{A}^T\mathbf{B}| \quad (26)$$

a useful special case is

$$|\mathbf{I}_N + \mathbf{a}\mathbf{b}^T| = 1 + \mathbf{a}^T\mathbf{b}$$

To prove the equation 26, we notice that

$$\begin{aligned} \begin{bmatrix} \mathbf{I}_N & \mathbf{0} \\ \mathbf{A}^T & \mathbf{I}_M \end{bmatrix} \begin{bmatrix} \mathbf{I}_N & \mathbf{B} \\ -\mathbf{A}^T & \mathbf{I}_M \end{bmatrix} &= \begin{bmatrix} \mathbf{I}_N & \mathbf{B} \\ \mathbf{0} & \mathbf{I}_M + \mathbf{A}^T\mathbf{B} \end{bmatrix} \\ \begin{bmatrix} \mathbf{I}_N & \mathbf{B} \\ -\mathbf{A}^T & \mathbf{I}_M \end{bmatrix} \begin{bmatrix} \mathbf{I}_N & \mathbf{0} \\ \mathbf{A}^T & \mathbf{I}_M \end{bmatrix} &= \begin{bmatrix} \mathbf{I}_N + \mathbf{B}\mathbf{A}^T & \mathbf{B} \\ \mathbf{0} & \mathbf{I}_M \end{bmatrix} \end{aligned}$$

So we have

$$\det \begin{bmatrix} \mathbf{I}_N & \mathbf{B} \\ \mathbf{0} & \mathbf{I}_M + \mathbf{A}^T\mathbf{B} \end{bmatrix} = \det \begin{bmatrix} \mathbf{I}_N + \mathbf{B}\mathbf{A}^T & \mathbf{B} \\ \mathbf{0} & \mathbf{I}_M \end{bmatrix}$$

which equals to

$$|\mathbf{I}_M + \mathbf{A}^T\mathbf{B}| = |\mathbf{I}_N + \mathbf{B}\mathbf{A}^T| = |\mathbf{I}_N + \mathbf{A}\mathbf{B}^T|$$

where we use the result $\det \mathbf{A} = \det \mathbf{A}^T$.