

Notes 3-Linear Models for Regression

Junda Du

2019-8-28

1 Linear Basis Function Models

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) \stackrel{\phi_0(\mathbf{x})=1}{=} \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}),$$

where $\mathbf{w} = (w_0, \dots, w_{M-1})^T$, $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$, $y(\mathbf{x}, \mathbf{w})$ is called linear models (linear for \mathbf{w}), but $\boldsymbol{\phi}(\mathbf{x})$ is the nonlinear function, playing a role in pre-processing and feature extraction.

The basis function $\phi_j(x)$ have many choice, for example

- $\phi_j(x) = x^j$, which is used for polynomial regression in Chapter 1. Its limitation is that they are global functions of the input variables, so that changes in one region of input space affect all other regions.??
- 'Gaussian' basis functions $\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$. μ_j govern the locations of the basis functions in input space, s governs their spatial scale, and the normalization coefficient is unimportant because of w_j .
- Sigmoidal basis function $\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$, where $\sigma(a) = \frac{1}{1 + \exp(-a)}$ is the logistic sigmoid function. Equivalently, we can use 'tanh' function $\tanh(a) = 2\sigma(a) - 1$
- Fourier basis: each function represents a specific frequency and has infinite spatial extent. It's interesting to consider basis functions that are localized in both space and frequency, leading to a class of functions known as *wavelets*.

1.1 Maximum likelihood and least squares

Maximum likelihood steps (least squares appears in step 3):

1. Target variable probability distribution given parameters:

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = N(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}), y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

2. Joint probability given a series of observations $\mathbf{X} = x_1, \dots, x_N$ with corresponding target values $\mathbf{t} = t_1, \dots, t_N$:

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

3. Take the logarithm of the likelihood function:

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \ln N(t_n|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}),$$

where $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$, which corresponds to least squares.

4. Solving for parameters:

- Derivate log likelihood function with respect to \mathbf{w} :

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n)^T = 0$$

Solving for \mathbf{w} we obtain: $\mathbf{w}_{ML} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t}$, which is known as the *normal equations* for the least squares problem.

$\boldsymbol{\Phi}$ is an $N \times M$ matrix, called *design matrix*, $\Phi_{nj} = \phi_j(\mathbf{x}_n)$. And the quantity $(\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T$ is known as the *Moore-Penrose pseudo-inverse* of the matrix $\boldsymbol{\Phi}$.

- Derivate $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n)\}^2$ with respect to w_0 equal to zero, and solving for w_0 , we obtain

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j, \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$$

- Derivate log likelihood function with respect to β , giving

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

1.2 Geometry of least squares

1.3 Sequential learning

1.4 Regularized least squares

Add a regularization term to an error function $E_W(w) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$ in order to control over-fitting, so

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

Accordingly, solving for \mathbf{w} as before, we obtain $\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$.

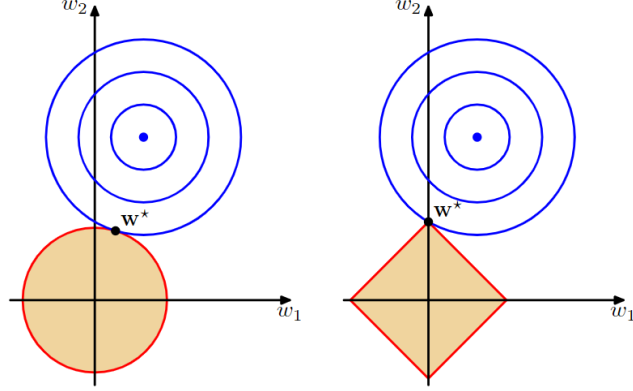
A more general regularization term is $E_W(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^M |w_j|^q$, so

$$E(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q.$$

And minimizing it is equivalent to minimizing the unregularized sum-of-squares error subject to the constraint $\sum_{j=1}^M |w_j|^q \leq \eta$, which can be found by Lagrange multipliers.

The case of $q = 1$ leads to a sparse model, i.e. some of the coefficients w_j are driven to zero. The origin can be seen from the following picture.

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}^* . The lasso gives a sparse solution in which $w_1^* = 0$.



1.5 Multiple outputs

If target vector \mathbf{t} is a K -dimensional column vector, the approach is to use the same set of basis function to model all of the components of the target vector so that $\mathbf{y}(\mathbf{x}, \mathbf{w}) = \mathbf{W}^T \phi(\mathbf{x})$, then the conditional distribution of the target vector is $p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = N(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1} \mathbf{I})$ assuming that is an isotropic Gaussian. The log likelihood function is

$$\ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) = \sum_{n=1}^K \ln N(\mathbf{t}_n|\mathbf{W}^T \phi(\mathbf{x}_n)) = \frac{NK}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n)\|^2$$

Solving for \mathbf{W} , giving $\mathbf{W}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$, examine this result for each target variable \mathbf{t}_k , we have $\mathbf{w}_k = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_k$, which leads to a decoupling into K independent regression problems when sharing the pseudo-inverse matrix.

2 The Bias-Variance Decomposition

It's unreasonable to minimize the regularized error function with respect to λ , a good approach is marginalizing over parameters in a Bayesian setting. Before doing so, it's instructive to consider a frequentist viewpoint of the model complexity issue, known as the *bias-variance*.

The expected squared loss can be written in the form

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

where $y(\mathbf{x})$ is which we seek a solution for to make first term a minimum, $h(\mathbf{x})$ is the optimal prediction given conditional distribution and the loss function, such as $h(\mathbf{x}) = E[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt$ given $p(t|\mathbf{x})$ and the squared loss function.

Supposed we have a large number of data sets each of size N , for any given data set D we can run learning algorithm and obtain a prediction function $y(\mathbf{x}; D)$. For a given data set D , the second term is about noise, independent of $y(\mathbf{x})$. Consider the integrand of the first term $\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2$, we introduce the quantity $E_D[y(\mathbf{x}; D)]$ because this term is dependent on the particular data set D . We obtain

$$\begin{aligned} & \{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)] + E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)]\}^2 + \{E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 \\ &+ 2\{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)]\}\{E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}, \\ &\text{then } E_D[\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2] = \underbrace{E_D[\{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)]\}^2]}_{\text{variance}} + \underbrace{E_D[\{E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2]}_{\text{bias}^2} \end{aligned}$$

$$E[L] = \text{bias}^2 + \text{variance} + \text{noise}, \text{ where } \text{bias}^2 = \int \{E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x},$$

$$\text{variance} = \int E_D[\{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)]\}^2] p(\mathbf{x}) d\mathbf{x}, \text{ noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

As we can see, there is a trade-off between bias and variance with flexible models (small λ) having high bias and low variance, and relatively rigid models (large λ) having low bias and high variance.

3 Bayesian Linear Regression

We turn to a Bayesian treatment of linear regression, which will avoid the over-fitting problem of maximum likelihood, and which will also lead to automatic methods of determining model complexity using the training data alone.

3.1 Parameter distribution

We want to obtain $p(\mathbf{w}|\mathbf{t})$. Noting the likelihood function $p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$

in 3.1.1, the corresponding conjugate prior is therefore given by the distribution of the form $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$ having the mean \mathbf{m}_0 and covariance \mathbf{S}_0 .

Making use of the result in 2.3.3, we can obtain

$$p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N), \mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t}), \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

For the remainder of this chapter, we shall consider a zero-mean isotropic Gaussian $p(\mathbf{w}|\alpha) = N(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$ to simplify the treatment. And corresponding posterior distribution over \mathbf{w} is given by

$$p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}, \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi.$$

The log posterior distribution is given by the sum of the log likelihood and log prior,

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const.}$$

Of course, it applies to sequential learning, i.e. it's easy to update the result when a new observation arrives.

There is a general Gaussian form of prior over the parameter \mathbf{w}

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

In our discussion of $q = 2$, the mode of posterior distribution is equal to mean.

3.2 Predictive distribution

We want to obtain $p(t|\mathbf{t})$. Given $p(t|\mathbf{w}, \beta)$ and $p(\mathbf{w}|\mathbf{t}, \alpha, \beta)$, the *predictive distribution* is defined by

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

Making use of the result in 2.3.3, it takes the form

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = N(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})), \sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}),$$

where $\sigma_N^2(\mathbf{x})$ is associated with data noise and parameters \mathbf{w} and they are independent so are additive. Also note that the level of uncertainty decreases as more data points are observed and $\sigma_{N+1}^2(\mathbf{x}) \leq \sigma_N^2(\mathbf{x})$. Problems:

- If we used localized basis functions such as Gaussians, then in region away from the basis function centre, the contribution from the second term in the predictive variance will go to zero. Thus, the model becomes very confident in its predictions when extrapolating outside the region occupied by the basis functions, which is generally an undesirable behaviour. This problem can be avoided by adopting an alternative Bayesian approach to regression known as a Gaussian process.
- If both \mathbf{w} and β are treated as unknown, then we can introduce a conjugate prior distribution $p(\mathbf{w}, \beta)$, which is given by Gaussian-gamma distribution discussed in 2.3.6.

3.3 Equivalent kernel

We want to introduce the equivalent kernel method by linear prediction. It's intuitive, jianjiede?? and have some properties.

We using $p(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$ to predict the target given the input, and we can estimate \mathbf{w} using \mathbf{m}_N in 3.3.1. Substitute it into the predictive function we can get

$$\begin{aligned}
 y(\mathbf{x}, \mathbf{m}_N) &= \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n \\
 &\stackrel{\substack{k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') \\ \text{equivalent kernel}}}{=} \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n, [\text{Property 1 : } \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1] \\
 &\stackrel{\substack{\psi(\mathbf{x}) = \beta^{1/2} \mathbf{S}_N^{1/2} \phi(\mathbf{x}) \\ \text{inner product transform}}}{=} \psi(\mathbf{x})^T \psi(\mathbf{x}), [\text{Property 2 : } k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})]
 \end{aligned}$$

We can obtain some insight from equivalent kernel:

- Regression functions, such as this, which makes predictions by taking combinations of the training set target values are known as *linear smoothers*. And the equivalent kernel is called smoother matrix.
- *Localization property??*: the weight of target is dependent on the defined distance[i.e. $\psi(\mathbf{x})$] between the data points and \mathbf{x} with closer data points has higher weight.
- The localization property holds not only for the localized Gaussian basis functions but also for the nonlocal polynomial and sigmoidal basis functions.??
- Instead of introducing a set of basis functions, which implicitly determines an equivalent kernel, we can instead define a localized kernel directly and use this to make predictions for new input given the observed training set.

4 Bayesian Model Comparison

There are many models to predict, we can choose the most accurate, known as *model selection*; we can also make use of the sum and product rules by

$$p(t|\mathbf{x}, D) = \sum_{i=1}^L p(t|\mathbf{x}, M_i, D) p(M_i|D), \text{ which is an example of a } \textit{mixture distribution} \text{ weighted by } p(M_i|D).$$

We need to obtain $p(M_i|D)$ by $p(M_i|D) \propto p(M_i)p(D|M_i)$, where $p(D|M_i)$

is important and known as *model evidence*. Then we try to obtain $p(D|M_i)$.

$$\begin{aligned}
p(D|M_i) &= \int p(D|\mathbf{w}, M_i) p(\mathbf{w}|M_i) d\mathbf{w} \\
\left[\text{Ignore } M_i \right] p(D) &= \int p(D|w) p(w) dw \\
\frac{\text{p(w) is peaked around } w_{MAP} \text{ with width } \Delta w_{prior}}{\text{p(D|w) is peaked around } w_{MAP} \text{ with width } \Delta w_{posterior}} p(D|w_{MAP}) \frac{\Delta w_{posterior}}{\Delta w_{prior}} \\
\left[\text{Taking logs} \right] \ln p(D) &= \ln p(D|w_{MAP}) + \ln \left(\frac{\Delta w_{posterior}}{\Delta w_{prior}} \right) \\
\left[\text{The model has } M \text{ parameters} \right] \ln p(D) &= \ln p(D|\mathbf{w}_{MAP}) + M \ln \left(\frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)
\end{aligned}$$

Here, we can see that the maximum evidence is given by the trade-off between these two competing terms. This is why Bayesian framework avoids the over-fitting problems.

5 The Evidence Approximation

In a fully Bayesian treatment, we don't select any model, but consider all models, i.e. integrate over all parameters. So

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

We have obtained $p(t|\mathbf{w}, \beta)$, $p(\mathbf{w}|\mathbf{t}, \alpha, \beta)$, so we need to obtain $p(\alpha, \beta|\mathbf{t})$ and $p(\alpha, \beta|\mathbf{t}) \propto p(\mathbf{t}|\alpha, \beta) p(\alpha, \beta)$ from Bayes' theorem. The important term is $p(\mathbf{t}|\alpha, \beta)$ because chances are that the prior is relatively flat or the training data set is large.

5.1 Evaluation of the evidence function

$$\begin{aligned}
p(\mathbf{t}|\alpha, \beta) &= \int p(\mathbf{t}|\mathbf{w}, \beta) p(\mathbf{w}|\alpha) d\mathbf{w} \\
\frac{3.1.1, 3.3.1, 2.3.3}{\left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2}} \int \exp\{-E(\mathbf{w})\} d\mathbf{w}
\end{aligned}$$

where

$$\begin{aligned}
E(\mathbf{w}) &= \beta E_D(\mathbf{w}) + \alpha E_W(\mathbf{w}) \\
&= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \\
\frac{\text{complete square}}{E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N)}, \\
\mathbf{A} &= \alpha \mathbf{I} + \beta \Phi^T \Phi = \nabla \nabla E(\mathbf{w}) = \mathbf{S}_N^{-1}, \\
E(\mathbf{m}_N) &= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N, \mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t}
\end{aligned}$$

So

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

Finally,

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

For the example in chapter 1, fix α , change M , the evidence favours the model with $M = 3$.

5.2 Maximizing the evidence function

First consider the maximization of $p(\mathbf{t}|\alpha, \beta)$ with respect to α . Obtaining $\frac{d}{d\alpha} \ln p(\mathbf{t}|\alpha, \beta)$, needs $\frac{d}{d\alpha} \ln |\mathbf{A}|$, then needs to define the following eigenvector equation $(\beta \mathbf{\Phi}^T \mathbf{\Phi}) \mathbf{u}_i = \lambda_i \mathbf{u}_i$, so \mathbf{A} has eigenvalues $\alpha + \lambda_i$, so

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

So

$$\begin{aligned}\frac{d}{d\alpha} \ln p(\mathbf{t}|\alpha, \beta) &= \frac{M}{2\alpha} - \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha} = 0 \\ \Leftrightarrow \alpha \mathbf{m}_N^T \mathbf{m}_N &= M - \alpha \sum_i \frac{1}{\lambda_i + \alpha} = \gamma, \\ \Leftrightarrow \alpha &= \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}, \gamma = M - \alpha \sum_i \frac{1}{\lambda_i + \alpha} = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}\end{aligned}$$

Similarly, we obtain

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

5.3 Effective number of parameters

γ represent the number of the effective number because for $\lambda \gg \alpha$, $\frac{\lambda_i}{\alpha + \lambda_i}$ is 1, which means λ_i is effective, for $\lambda \ll \alpha$, $\frac{\lambda_i}{\alpha + \lambda_i}$ is 0, which means λ_i is non-effective.

For the case of $N \gg M$, $\gamma = M$ because $\mathbf{\Phi}^T \mathbf{\Phi}$ involves an implicit sum over data points and so the eigenvalues λ_i increase with the size of the data set. So the re-estimation equations for α, β becomes

$$\alpha = \frac{M}{2E_W(\mathbf{m}_N)}, \beta = \frac{N}{2E_D(\mathbf{m}_N)}$$

6 Limitations of Fixed Basis Functions

The basis functions $\phi_j(\mathbf{x})$ are fixed before the training data set is observed and is a manifestation of the curse of dimensionality. As a consequence, the number of basis functions needs to grow rapidly.

Fortunately, there are two properties of real data sets that we can exploit to help alleviate this problem. First, the real dimension of \mathbf{x} is smaller than input space because of the strong correlations between the input variables. The second property is that the target may have significant dependence on only a small number of possible directions.