

# Linear Models for Regression

A self-study materials for PRML [1]

Jisung Lim<sup>1</sup>

<sup>1</sup>B.S. Candidate of Industrial Engineering  
Yonsei University, South Korea.

8th February, 2017



YONSEI  
UNIVERSITY



Christopher M. Bishop. **Pattern Recognition and Machine Learning**. Springer, 2006.

# Summary

## 1 Introduction

- Main concept of linear regression
- Linear models and regression

## 2 Frequentist - linear regression

- Prediction by Maximum Likelihood Estimation
- Geometrical Interpretation of Least Squares
- Sequential learning
- Regularized least squares
- Multiple outputs

## 3 Frequentist - model complexity

- A frequentist viewpoint of the model complexity
- A bias-variance tradeoffs

## 4 Bayesian - linear regression

- Bayesian approach in linear regression
- Bayesian treatment
- Equivalent kernel

## 5 Bayesian - model complexity

- A Bayesian viewpoint of the model complexity
- The evidence approximation
- An elegant interpretation of  $\alpha$  and  $\gamma$



# Main Concept of Linear Regression

## The main concept of linear regression

The goal of regression is to predict the value of one or more continuous target variable  $t$  (or  $\mathbf{t}$  for more than one variable) given the value of a  $D$ -dimensional vector  $\mathbf{x}$  of input variables. Among them, linear regression model is a broad class of functions which share the property of being linear functions with respect to the adjustable parameters  $\theta_j$ .



# Linear Models and Regression

## Three components of linear regression

We may see the linear regression as a process comprise of three components

### 1 N observations

At first, We retrieve a training data set of  $N$  observations  $\mathbf{x}_n$  together with corresponding target values  $t_n$ .

$$\mathcal{D} = \{(t_1, \mathbf{x}_1), \dots, (t_N, \mathbf{x}_N)\} \quad (1)$$

### 2 Linear Models for Regression

Given  $N$  data set, we will construct linear regression models that comprise a broad class of functions which share the property of being linear functions with respect to the adjustable parameters  $\theta_j$ . We may extend the class of models by using basis functions  $\phi_j(\cdot)$ .

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

### 3 Prediction

Then, we can construct a function  $y(\mathbf{x})$  which constitute the predictions for the target value  $t$ .

$$\hat{t} = y(\mathbf{x}) \quad (3)$$

or more generally, from a probabilistic perspective we can construct probability distribution of  $t$  given  $\mathbf{x}$

$$p(t|\mathbf{x}) \quad (4)$$

which expresses our uncertainty about the value of  $t$  for each value of  $\mathbf{x}$ .

# The Form of Linear Models

Linear regression model is a broad class of functions which share the property of being linear functions of the adjustable parameters  $\theta_j$ . The form of linear regression models varies from the form of linear combination of parameters  $\theta_j$  and input variables  $x_j$  to the form of linear combination of the parameters and a fixed set of nonlinear functions.

Given  $\mathbf{x} = (x_1, \dots, x_D)^T$ , the simplest linear model for regression forms:

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \sum_{j=1}^D \theta_j x_j \quad (5)$$

which is often simply known as *linear regression*. When augmented with other forms of basis function expansion, it can model also nonlinear relationships.

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \sum_{j=1}^{M-1} \theta_k \phi_j(\mathbf{x}) \quad (6)$$

or equivalently, introducing dummy basis function  $\phi_0(\mathbf{x}) = 1$  so that

$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) \quad (7)$$

the fixed set of nonlinear functions of the input variables  $\boldsymbol{\phi}(\mathbf{x})$  is known as *basis functions* and  $\theta_0$  as *bias parameter* which allows for any fixed offset in the data.

# Linear Models and Basis Function

In machine learning approach, we will apply some form of fixed pre-processing or feature extraction, to the original data variables.

$$\begin{aligned} \text{original variables : } \mathbf{x} &= (x_1, \dots, x_N)^T \\ \text{features : } \phi(\mathbf{x}) &= (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T \end{aligned} \quad (8)$$

## Choose of basis function

- Polynomial basis function (ex)  $\phi_j(\mathbf{x}) = x^j$   
Limitation of polynomial basis function is that they are global functions of the input variables.
- Spline function (ex)  $\phi_j(\mathbf{x}) = f_i(\mathbf{x})$  for  $\mathbf{x} \in \mathcal{R}_i$   
This can be resolved by dividing the input space up into regions and fit a different polynomial  $f_i(\mathbf{x})$  in each region  $\mathcal{R}_i$ .
- Gaussian basis function (ex)  $\phi_j(x) = \exp\{-(x - \mu^2)^2/2s^2\}$   
A probabilistic interpretations are not required, and in particular the normalization coefficient is unimportant because of the parameters  $w_j$  which is to be multiplied to each basis function.
- Sigmoid and tanh basis function (ex)  $\phi_j(x) = \sigma((x - \mu_j)/s)$  where  $\sigma(a)$  is the logistic sigmoid function. Since  $\tanh(a) = 2\sigma(a) - 1$ , generalized linear combination of two basis functions can be identically expressed.
- Fourier basis function (ex) wavelets.



# Gaussian Noise Assumption

## setting 1. Additive gaussian noise

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

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(\epsilon \mid 0, \beta^{-1}) \quad (9)$$

Then, we may express the distribution of target value  $t$  as follows

$$p(t \mid \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(\epsilon \mid y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (10)$$

## setting 2. Squared loss function

If we set the loss function as a squared loss function as follows

$$L(t, y(\mathbf{x}, \mathbf{w})) = \{t - y(\mathbf{x}, \mathbf{w})\}^2 \quad (11)$$

and also we can define the expected loss as follows

$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x}, \mathbf{w})) p(t, \mathbf{x}) dt d\mathbf{x} = \iint \{t - y(\mathbf{x}, \mathbf{w})\}^2 p(t, \mathbf{x}) dt d\mathbf{x} \quad (12)$$

Now, to find optimal prediction  $y$ , set variational derivative of  $\mathbb{E}[L]$  with respect to  $y(\mathbf{x}, \mathbf{w})$  equals to zero  $\delta \mathbb{E}[L] / \delta y = 0$ , then we get  $y^* = \mathbb{E}[t \mid \mathbf{x}]$  which will be simply

$$y^* = \mathbb{E}[t \mid \mathbf{x}] = \int t p(t \mid \mathbf{x}) dt = y(\mathbf{x}, \mathbf{w}) \quad (13)$$





# Maximum Likelihood with Basis Function

Now, we will construct predictive function  $y(t|\mathbf{x}, \mathbf{w}, \beta)$  by using maximum likelihood estimation.

## $N$ Observations

At first, Lets consider the input data of  $N$  observations,

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \quad \text{and} \quad \mathbf{t} = (t_1, \dots, t_N)^T \quad (14)$$

## Find likelihood function

With i.i.d. assumption, we can get the following likelihood function,

$$L(\mathbf{w}) = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (15)$$

taking logarithm of both sides,

$$\begin{aligned} \ln L(\mathbf{w}) &= \sum_{n=1}^N \ln[\mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})] \\ &= \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi) - \frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \\ &= \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \end{aligned} \quad (16)$$

where  $E_D(\mathbf{w})$  is the sum-of-squares error function.



## Maximum likelihood with respect to $w$

### Maximize $L(w)$ with respect to $w$

As we've seen before, maximizing the likelihood function  $L(w)$  is equivalent to minimizing a sum-of-squares error function  $E_D(w)$  from (16), given the following assumptions

$$t \sim \mathcal{N}(t | y(x, w), \beta^{-1}) \quad \text{with i.i.d. condition} \quad (17)$$

where  $y(x, w) = w^T \phi(x_n)$ . Or we may find same result with the gradient of the log likelihood function,

$$\nabla p(\mathbf{t} | w, \beta) = \sum_{n=1}^N [\{t_n - w^T \phi(x_n)\} \phi(x_n)^T] \quad (18)$$

Setting the gradient to zero gives

$$0 = \sum_{n=1}^N t_n \phi(x_n)^T - w^T \left( \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \right) \quad (19)$$

Solving for  $w$ , we obtain

$$w_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (20)$$



## Maximum likelihood with respect to $\beta$

### Maximize $L(\mathbf{w})$ with respect to $\beta$

We can also maximize the log likelihood function with respect to the noise precision parameter  $\beta$ , setting derivate of likelihood function with respect to  $\beta$  equal to zero as follows

$$\frac{\partial}{\partial \beta} [p(\mathbf{t} | \mathbf{w}, \beta)] = \frac{N}{2\beta} - E_D(\mathbf{w}) = 0 \quad (21)$$

Then we get,

$$\frac{1}{\beta_{\text{ML}}} = \frac{2}{N} E_D(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_n)\}^2 \quad (22)$$

### Predictive Function

Hence, we get predictive function which is given by

$$p(t | \mathbf{x}, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \left( \frac{\beta_{\text{ML}}}{2\pi} \right)^{1/2} \exp \left[ -\frac{\beta_{\text{ML}}}{2} \left\{ t - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}) \right\}^2 \right] \quad (23)$$



# Normal Equation for the least squares problem

## Normal Equation

Now consider some conventional terms with the solution (20) for  $w$

$$w_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

which are known as the *normal equations* for the least squares problem.

## Design Matrix

Here  $\Phi$  is an  $N \times M$  matrix, called the *design matrix*,

$$\Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{bmatrix} = [\varphi_0 \quad \cdots \quad \varphi_{M-1}] \quad (24)$$

where

$$\begin{aligned} \phi(\mathbf{x}_n) &= (\phi_0(\mathbf{x}_n), \dots, \phi_{M-1}(\mathbf{x}_n))^T \\ \varphi_j &= (\phi_j(\mathbf{x}_1), \dots, \phi_j(\mathbf{x}_N))^T \end{aligned} \quad (25)$$

## Moore-Penrose pseudo-inverse

And then, the quantity

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T \quad (26)$$

is known as the *Moore-Penrose pseudo-inverse* of the matrix  $\Phi$ .

## The Role of Bias Parameter $w_0$

We can get some insight of the role of *Bias Parameter*  $w_0$  by solving the equation

$$\frac{d}{dw_0} E_D(\mathbf{w}) = 0 \quad (27)$$

At first, let us make the bias parameter to be explicit.

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right\}^2 \quad (28)$$

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

$$\begin{aligned} w_0 &= \frac{1}{N} \sum_{n=1}^N t_n - \sum_{j=1}^{M-1} \left[ w_j \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n) \right] \\ &= \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \end{aligned} \quad (29)$$

Thus, the bias  $w_0$ , as a kind of offset, compensates for the gap between the averages of the target values  $\bar{t}$  and the linear combination of adaptive parameters  $w_j$  and the average of the basis functions  $\bar{\phi}_j$ .

# Geometry of least squares

## Geometry of least squares

Let us consider an  $N$ -dimensional space whose axes are given by the  $t_n$ , so that  $\mathbf{t} = (t_1, \dots, t_N)^T$  is a specific vector in this space.

$$\mathbf{t} = (t_1, \dots, t_N)^T \in \mathbb{R}^N \quad (30)$$

Each basis function  $\phi_j(x_n)$  may be represented as a vector  $\varphi_j$  in the same space, when evaluated at the  $N$  data points.

$$\varphi_j \in \mathbb{R}^N \quad \text{where} \quad \Phi = [\varphi_0 \quad \cdots \quad \varphi_{M-1}] \quad (31)$$

Since  $\mathbf{y}$  is an arbitrary linear combination of the vector  $\varphi_j$ , it lies on the  $M$ -dimensional subspace  $\mathcal{S}$ .

$$\mathbf{y} = \Phi \mathbf{w} = [\varphi_0 \quad \cdots \quad \varphi_{M-1}] \begin{bmatrix} w_0 \\ \vdots \\ w_{M-1} \end{bmatrix} = \sum_{j=0}^{M-1} w_j \varphi_j \quad (32)$$

Hence, intuitively, the sum-of-squares error may be understood as the squared Euclidean distance between  $\mathbf{y}$  and  $\mathbf{t}$ . Thus the least-squares solution for  $\mathbf{w}$  corresponds to the choice of  $\mathbf{y}$  that lies in subspace  $\mathcal{S}$  and that is closest to  $\mathbf{t}$ .

# Sequential learning

In cases where the training data set is very large or data is received in a stream, a direct solution using the normal equations may not be possible. An alternative approach is to use sequential algorithms, also known as on-line algorithms. Sequential learning is appropriate for realtime applications in which the data observations are arriving in a continuous stream, and predictions must be made before all of the data points are seen. At first, we will consider the data points  $D_n = (n = 1, 2, 3, \dots)$  one at a time and then update model parameters after each such consideration.

$$\mathcal{D} = \{D_1, \dots, D_n, \dots\} = \{(t_1, \mathbf{x}_1), \dots, (t_n, \mathbf{x}_n), \dots\} \quad (33)$$

After presentation of  $D_n$ , the stochastic gradient descent algorithm updates the parameter vector  $\mathbf{w}$  using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n \quad (34)$$

The total error function  $E$  is the sum of a given error function  $E_n$ , evaluated from the  $n$ -th data point  $D_n$ .

$$E = \sum_{\forall n} E_n \quad \text{where} \quad E_n = d(t_n, y(\mathbf{x}_n)) = \{t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n)\}^2 \quad (35)$$

and the algorithm is given as follows, where  $\eta$  represents a learning rate.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \{t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n) \quad (36)$$

This is known as *least-mean-squares* or the *LMS algorithm*.



# Regularized Least Squares

Let us adding a regularization term to an error function in order to prevent model from over-fitted. The total error function to be minimized takes the form

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \quad (37)$$

where  $\lambda$  is the regulariation coefficient that controls the relative importance between the data-dependent error  $E_D(\mathbf{w})$  and the regularization term  $E_W(\mathbf{w})$ .

One of the simple forms of regularizer is given by the sum-of-squares of the weight vector elements

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}. \quad (38)$$

If we consider the sum-of-squares error function given by

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \quad (39)$$

then the total error function becomes

$$E(\mathbf{w}) = E_D(\mathbf{w}) + E_W(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}. \quad (40)$$





## Other choices of regularizer

### Generalized regularizer

A more general regularizer is sometimes used, for which the regularized error takes the form

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q \quad (41)$$

where  $q = 2$  corresponds to the quadratic regularizer. There are conventional terms for linear and quadratic regularizers, respectively, *lasso* and *ridge*.

### Minimizing regularized model and constrained optimization

For sufficiently large  $\lambda = \lambda^*$ , we have to minimize regularized model (38) to obtain the solution  $\mathbf{w}^*$ .

$$\underset{\mathbf{w}_j}{\text{Minimize}} \quad \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \lambda^* \sum_{j=1}^M |w_j|^q \quad (42)$$

This optimization problem can be seen as a Lagrange multipliers, so that it can be represented by constrained optimization problem.

$$\begin{aligned} &\underset{\mathbf{w}_j}{\text{Minimize}} \quad \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \\ &\text{subject to} \quad \sum_{j=1}^M |w_j|^q \leq \eta^* \end{aligned} \quad (43)$$



# Multiple Outputs

# A Frequentist Viewpoint of The Model Complexity

## Model Complexity

The use of MLE can lead to severe over-fitting if complex models are trained using data sets of limited size. For releasing the over-fitting problem, we've discussed two approaches:

- limiting the number of basis functions
- introducing regularization terms

But, *the first approach* has side effect of limiting the flexibility of the model to capture interesting and important trends in the data. Although, *the second approach* reveals the problem of choosing the appropriate coefficient  $\lambda$ .

## A frequentist viewpoint of the model complexity

Now we consider the over-fitting issue and model complexity. As we've seen before, an over-fitting issue doesn't incur when we marginalize over parameters in a Bayesian setting. In the frequentist setting, however, they handle the model complexity issue with a special viewpoint, known as the *bias-variance* trade-off. At first, we will consider the frequentist view point.



## Loss Function for Regression

In decision theory, we consider the family of functions  $y(\mathbf{x})$  which is estimate of the value of  $t$ , for each input  $\mathbf{x}$ . Generally, the optimal function  $y^*(\mathbf{x})$  can be chosen by minimizing the loss function. And a popular choice of the loss function is the squared loss function which is given by

$$L(t, y(\mathbf{x})) = \{t - y(\mathbf{x})\}^2 \quad (44)$$

and the expected loss function is given by

$$\begin{aligned} \mathbb{E}[L] &= \iint L(t, y(\mathbf{x})) p(t, \mathbf{x}) d\mathbf{x} dt \\ &= \iint \{t - y(\mathbf{x})\}^2 p(t, \mathbf{x}) d\mathbf{x} dt \end{aligned} \quad (45)$$

Now we want to find the optimal function  $y^*(\mathbf{x})$ , which minimizes the expected loss (42), and we can get by setting the variational derivative  $\delta \mathbb{E}[L] / \delta y(\mathbf{x})$  equals to zero as follows

$$\begin{aligned} \frac{\delta \mathbb{E}[L]}{\delta y(\mathbf{x})} &= 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) dt = 0 \\ \iff y(\mathbf{x}) \int p(\mathbf{x}, t) dt &= \int t p(\mathbf{x}, t) dt \\ \iff y(\mathbf{x}) = \int t \frac{p(\mathbf{x}, t)}{p(\mathbf{x})} dt &= \int t p(t|\mathbf{x}) dt = \mathbb{E}[t|\mathbf{x}] \end{aligned} \quad (46)$$



## Loss Function for Regression (conti. . .)

From the result of (43), for convenience, denote the conditional expectation by  $h(\mathbf{x})$ , which is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x})dt \quad (47)$$

Now the expected loss function can be written in the form

$$\begin{aligned} \mathbb{E}[L] &= \mathbb{E}[\{t - h(\mathbf{x})\}^2] + \mathbb{E}[\{h(\mathbf{x}) - y(\mathbf{x})\}^2] + \mathbb{E}[2\{t - h(\mathbf{x})\}\{h(\mathbf{x}) - y(\mathbf{x})\}] \\ &= \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \end{aligned} \quad (48)$$

since

$$\begin{aligned} \mathbb{E}[\{t - h(\mathbf{x})\}\{h(\mathbf{x}) - y(\mathbf{x})\}] &= \int \{y(\mathbf{x}) - h(\mathbf{x})\} \int \{h(\mathbf{x}) - t\} p(\mathbf{x}, t) dt d\mathbf{x} \\ &= \int \{y(\mathbf{x}) - h(\mathbf{x})\} \left\{ h(\mathbf{x}) p(\mathbf{x}) - \int tp(\mathbf{x}, t) dt \right\} d\mathbf{x} \\ &= 0 \quad \left( \because h(\mathbf{x}) = \int tp(t|\mathbf{x}) dt \right) \end{aligned}$$



## Loss Function for Regression (conti. . .)

Now consider the equation (45)

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

**The second term:** intrinsic, invariant noise

$$\mathbb{E}[\{h(\mathbf{x}) - t\}^2] = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- The second term is independent of  $y(\mathbf{x})$ .
- This quantity arises from the intrinsic noise on the data.
- It represents the minimum achievable value of the expected loss.

**The first term:** Error arose from the choice for the function  $y(\mathbf{x})$

$$\mathbb{E}[\{y(\mathbf{x}) - h(\mathbf{x})\}^2] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

- The first term depends on our choice for the function  $y(\mathbf{x})$ .
- We will seek a solution for  $y(\mathbf{x})$  which makes this term a minimum.
- Because it is nonnegative, the smallest that we can hope to make this term is zero.
- If the data are sufficient enough, we could find the regression function  $h(\mathbf{x})$  in any desired degree of accuracy.

## Frequentist treatment for modeling the $h(\mathbf{x})$

In a viewpoint of frequentist, when we model the function  $h(\mathbf{x})$  using the parametric function  $y(\mathbf{x}, \mathbf{w})$ , since the frequentist treatment involves making a point estimate of  $\mathbf{w}$  based on the data set  $\mathcal{D}$ , they try to understand the uncertainty of the estimate as the frequency when we suppose a large number of data sets each of which satisfies i.i.d. condition.

Let us consider the case where the size of data set is  $N$ , the number of data set is  $L$ , and each data set is drawn from the distribution  $p(t, \mathbf{x})$  independently.

$$\mathcal{D}^{(l)} = \{(t_1, \mathbf{x}_1), \dots, (t_N, \mathbf{x}_N)\} \quad \text{where } l = 1, \dots, L \quad \text{and} \quad \mathcal{D}^{(l)} \sim p(t, \mathbf{x}) \quad (50)$$

Then we could run the learning algorithm to obtain a prediction function  $y(\mathbf{x}; \mathcal{D})$  and then, from the first term of the equation (45), the integrand will be given by

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 = \{(y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]) + (\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}))\}^2 \quad (51)$$

for any given specific data set  $\mathcal{D}$ .

## Frequentist treatment for modeling the $h(\mathbf{x})$ (conti. . .)

Since the quantity (48) will be dependent on the particular data set  $\mathcal{D}$ , we should take its average over all data sets.

$$\begin{aligned}\mathbb{E}[\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] &= \mathbb{E} \left[ \{(y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]) + (\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}))\}^2 \right] \\ &= \{\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 + \mathbb{E} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]\}^2 \right]\end{aligned}\quad (52)$$

since

$$\mathbb{E}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]\} \{\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}] = 0 \quad (53)$$

Hence, the average quantity over all data sets is given by the form

$$\mathbb{E}[\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] = \underbrace{\{\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]\}^2 \right]}_{\text{variance}} \quad (54)$$

where the equation might be interpreted as the **bias-variance tradeoffs**.





## Frequentist treatment for modeling the $h(\mathbf{x})$ (conti. . .)

Hence, now we can consider the expected loss under the frequentist setting as follows

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise} \quad (55)$$

$$\text{where } (\text{bias})^2 = \int \{\mathbb{E}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y(\mathbf{x}; \mathcal{D})]\}^2 p(\mathbf{x}) d\mathbf{x} \quad (56)$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(t, \mathbf{x}) d\mathbf{x} dt$$

Since we only have finite data sets where  $l = 1, \dots, L$ , the bias and variance terms are given by

$$\begin{aligned} (\text{bias})^2 &= \frac{1}{N} \sum_{n=1}^N \{\bar{y}(\mathbf{x}_n) - h(\mathbf{x}_n)\}^2 \\ \text{variance} &= \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(\mathbf{x}_n) - \bar{y}(\mathbf{x}_n)\}^2 \end{aligned} \quad (57)$$

$$\text{where } \bar{y}(\mathbf{x}) = \mathbb{E}[y(\mathbf{x}; \mathcal{D})] = \frac{1}{L} \sum_{l=1}^L y^{(l)}(\mathbf{x})$$



# Frequentist treatment for modeling the $h(x)$ (conti. . .)

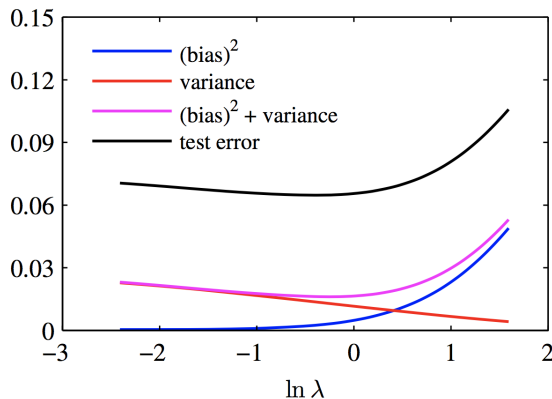


Figure 1: Loss function for the regularized least-square model.



# Bayesian Approach in Linear Regression

## What we've discussed so far

- Effective model complexity; governed by the number of basis function.
- Regularizer; governed by the value of the regularization coefficient.
- Validation set; too much expensive and wasteful.

## We, therefore, turn to a Bayesian treatment of linear regression, which

- will avoid the over-fitting problem of Maximum Likelihood approach.
- will determine the model complexity automatically, using the training data alone.

## What we will do

- Formulate our knowledge about the model itself, and prior belief for the uncertainty of parameters.

$$\text{Model : } y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}) \quad \text{Parameter : } p(\mathbf{w}|\alpha) \quad (58)$$

- Observe data of  $N$  observations

$$\mathcal{D} = \{(t_1, \mathbf{x}_1), \dots, (t_N, \mathbf{x}_N)\} \quad (59)$$

- Evaluate posterior distribution for the parameters

$$p(\mathbf{w}|\mathcal{D}) \propto L(\mathbf{w})p(\mathbf{w}) \quad (60)$$

- Make prediction

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta)p(\mathbf{w}|\mathbf{t}, \alpha, \beta)d\mathbf{w} \quad (61)$$



# Bayesian Treatment

First, assume we get the likelihood function from the data which has the form

$$\text{likelihood : } p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \sim \text{Gaussian distribution} \quad (62)$$

Since the likelihood function is given by a Gaussian distribution so that the conjugate prior and corresponding posterior is also given by a Gaussian distribution.

$$\text{prior : } p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad (63)$$

$$\text{posterior : } p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad (64)$$

where

$$\begin{aligned} \mathbf{m}_N &= \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^T \Phi \end{aligned} \quad (65)$$

Note that,

- $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$ , since the Gaussian is unimodal distribution whose mode is coincides with its mean.
- Also if we chose an infinitely broad prior  $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$  with  $\alpha \rightarrow 0$ , the posterior mean  $\mathbf{m}_N$  reduces to the value  $\mathbf{w}_{\text{ML}}$ .
- Similarly, if  $N = 0$ , the posterior and prior are same.
- Furthermore, from the concept of on-line learning, the posterior at any stage would become the prior at the next stage.

# Bayesian Treatment

## Set the prior

For simplicity, we will set the prior distribution as follows

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (66)$$

## Evaluate the posterior

We, therefore, get the mean and variance of the posterior given by

$$\begin{aligned} \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi \end{aligned} \quad (67)$$

## Maximum a posteriori approach

Since the posterior is proportional to the product of the prior and the likelihood function, the log of the posterior is given by

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const.} \quad (68)$$

where maximizing the posterior distribution with respect to  $\mathbf{w}$  is equivalent to the minimizing the regularized sum-of-squares error function.

# Bayesian Treatment

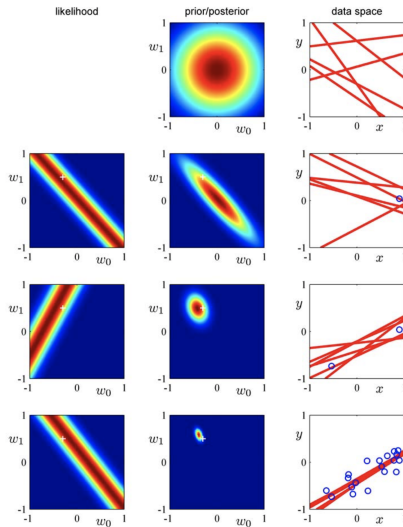


Figure 2: Illustration of the sequential Bayesian treatment.

# Predictive distribution

Indeed, marginalizing multiple solutions with respect to the distribution of the parameter  $\mathbf{w}$  lies at the heart of a Bayesian approach. We, therefore, get the predictive distribution given by

$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t, \mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \int p(t|\mathbf{w}, \mathbf{t}, \alpha, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \end{aligned} \quad (69)$$

where the function is not relevant to the parameter  $\mathbf{w}$ . Since the (62) involves the convolution of two Gaussian distributions, the predictive distribution could be the form

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \Phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (70)$$

where

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi^T \mathbf{S}_N \phi(\mathbf{x}) \quad (71)$$

Note that, the first term represents the noise on the data, whereas the second term reflects the uncertainty associated with the parameters  $\mathbf{w}$ . And also, if the size of data set becomes larger, the second term goes to zero.

# Predictive distribution

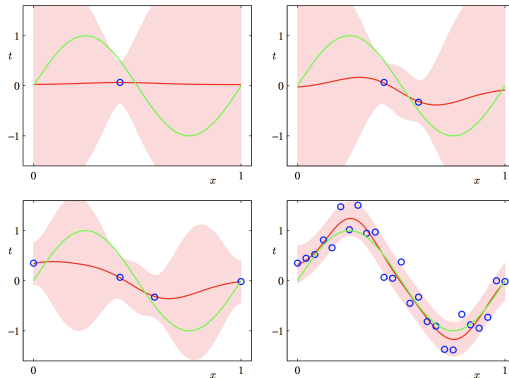


Figure 3: Examples of predictive distributions.



# Model comparison

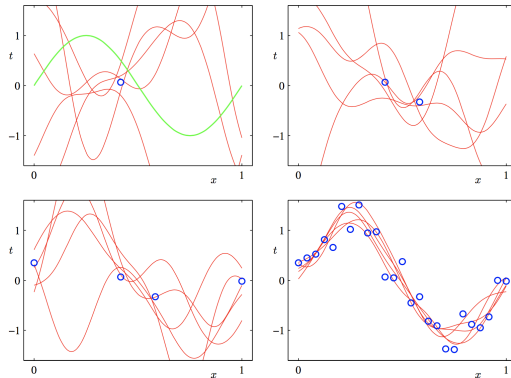


Figure 4: Plots of the function  $y(x, \mathbf{w})$  from the posterior distributions.

# Equivalent Kernel

We can see that the predictive mean as a specific model given by

$$y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n \quad (72)$$

where  $\mathbf{w} = \mathbf{m}_N$ . Then, the equation (65) may be interpreted as a linear combination of the target variables  $t_n$  from training set, so that we can write

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n \quad (73)$$

where the function

$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') \quad (74)$$

is known as *the smoother matrix* or *the equivalent kernel*. And, also *the equivalent kernel* could be obtained by

$$\begin{aligned} \text{cov}[y(\mathbf{x}), y(\mathbf{x}')] &= \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] \\ &= \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') \\ &= \beta^{-1} k(\mathbf{x}, \mathbf{x}') \end{aligned} \quad (75)$$

Note that, the predictive mean at nearby points will be highly correlated, whereas for more distant pairs of points the correlation will be smaller.

# A Bayesian Viewpoint of The Model Complexity

## What is model selection?

- Choosing the value of the hyperparameters.
- Choosing a model between alternative models.

## How?

- Frequentist validates their model selection using validation set in their Variance-bias framework.
- Bayesian use the probabilistic framework to express the uncertainty in the choice of model.



# Bayesian Model Comparison

## Uncertainty of the model

Now, suppose we wish to compare a set of  $L$  models  $\{\mathcal{M}_i\}$  and suppose that the data  $\mathcal{D}$  is generated from one of these models, but we are uncertain which one.

$$data : \mathcal{D} \quad model : \{\mathcal{M}_i\} \forall i = \{1, \dots, L\} \quad (76)$$

## Prior distribution for models

Each model  $\{\mathcal{M}_i\}$  refers to a probability distribution over the observed data  $\mathcal{D}$ . And our uncertainty over models would be expressed through a prior probability distribution

$$p(\mathcal{M}_i) \quad (77)$$

## Posterior distribution for models

Given a training set  $\mathcal{D}$ , we then wish to evaluate the posterior distribution which is calculated by the Bayes' theorem given by

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i) \quad (78)$$

The evidence term  $p(\mathcal{D}|\mathcal{M}_i)$  also called the marginal likelihood function because the parameters  $w$  are marginalized out.

# Predictive distribution

## Predictive distribution

Once we know the posterior distribution over models, the predictive distribution is given by

$$p(t|\mathbf{x}, \mathcal{D}) = \sum_{i=1}^L p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) \quad (79)$$

The predictive distribution is marginalized over the space of models.

- $p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})$ : Predictive distributions of individual models
- $p(\mathcal{M}_i|\mathcal{D})$ : Posterior probabilities of those models.

## Model selection

A simple approximation to the predictive distribution marginalized over models is to use the single, most probable model alone to make predictions. This is known as *model selection*.



# Evidence (Marginal Likelihood)

## Evidence (Marginal Likelihood)

For a model governed by a set of parameters  $\mathbf{w}$ , the model evidence is given by

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)d\mathbf{w} \quad (80)$$

- $p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)$ : probability of generating the data set  $\mathcal{D}$  from a model  $\mathcal{M}_i$
- $p(\mathbf{w}|\mathcal{M}_i)$ : prior distribution for parameters  $\mathbf{w}$

## Simple Approximation to Evidence

Now, to get some insight into the model evidence, we approximate the integral. At first, let us assume some simplified situation as follows

- The posterior is picked around the most probable value  $w_{\text{MAP}}$ , with width  $\Delta w_{\text{posterior}}$
- The prior is flat with width  $\Delta w_{\text{prior}}$

Then, we get

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)d\mathbf{w} \simeq p(\mathcal{D}|w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \quad (81)$$

## Evidence (Marginal Likelihood)

and so taking logs we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\mathbf{w}_{\text{WAP}}) + \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right) \quad (82)$$

With  $M$  parameters, all assumed to have the same  $\Delta w_{\text{posterior}}/\Delta w_{\text{prior}}$  ratio, we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\mathbf{w}_{\text{WAP}}) + M \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right) \quad (83)$$

- **Fitness:**  $\ln p(\mathcal{D}|\mathbf{w}_{\text{WAP}})$

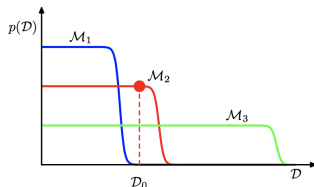
The first term represents the fit to the data given by the most probable parameter values. As we increase the complexity of the model, it increases because a more complex model is better able to fit the data.

- **Penalty:**  $M \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$

The second term represents how much the parameters tuned to the data. It penalizes the model according to its complexity. Since it has negative values, it decreases as we increase the complexity of the model.



# Bayesian Model Comparison



How a specific data set  $\mathcal{D}$  and its marginal likelihood  $p(\mathcal{D})$  can favour models of intermediate complexity? Let's consider the case where we observe a specific data  $\mathcal{D}_0$ , then we may assume three models as follows

Figure 5: Distribution of possible data set of models.

## ■ Low complexity $\mathcal{M}_1$

A simple model cannot fit the data well, whereas it assigns larger probability than more complex one.

## ■ Intermediate complexity $\mathcal{M}_2$

For specific data  $\mathcal{D}_0$  will have the largest probability when we choose a model that has intermediate complexity. Simpler one can not fit to the data, whereas more complex one can fit to the data but it has smaller probability.

## ■ High complexity $\mathcal{M}_3$

A complex model spreads its predictive probability, so that it assigns relatively small probability.



# Model Comparison and Bayes Factor

Let consider two models  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , where the correct answer is the model  $\mathcal{M}_1$ . For a given finite data set, it is possible that the Bayes factor flavour incorrect answer which is given by

$$\frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)} < 1 \quad (84)$$

However, if we average the Bayes factor over the distrubtion of data sets from true model

$$\int p(\mathcal{D}|\mathcal{M}_1) \ln \left( \frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)} \right) d\mathcal{D} = D_{\text{KL}}(\mathcal{M}_1||\mathcal{M}_2) \geq 0 \quad (85)$$

Thus on average the Bayes factor will always favour the correct model. For a given



# The Evidence Approximation

## The fully Bayesian approach

In the fully Bayesian approach, we would also specify a prior distribution over the hyperparameters, such as  $\alpha$  and  $\beta$ , so that they also can be marginalized in the Bayesian framework.

$$p(t|\mathbf{x}, \mathbf{t}) = \iiint p(t|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) p(\alpha, \beta|\mathbf{t}) d\mathbf{w} d\alpha d\beta \quad (86)$$

- **Our model:**  $p(t|\mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$
- **posterior over parameters:**  $p(\mathbf{w}|\mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$
- **posterior over hyperparameters:**  $p(\alpha, \beta|\mathbf{t})$

## Intractable solution and Evidence approximation

But the solution above is analytically intractable, so we need to evaluate the predictive model by approximating the model distribution. It is also called the *evidence approximation*

# The Evidence Approximation

## Assumptions for approximation

Let  $p(\alpha, \beta|\mathbf{t})$  is sharply peaked around values  $\hat{\alpha}$  and  $\hat{\beta}$ , then we get

$$p(t|\mathbf{t}) \simeq p(t|\mathbf{t}, \hat{\alpha}, \hat{\beta}) = \int p(t|\mathbf{w}, \hat{\beta})p(\mathbf{w}|\mathbf{t}, \hat{\alpha}, \hat{\beta})d\mathbf{w} \quad (87)$$

Now, assuming that the prior is relatively flat, from the Bayes' theorem

$$p(\alpha, \beta|\mathbf{t}) \propto p(\mathbf{t}|\alpha, \beta)p(\alpha, \beta) \quad (88)$$

## Evidence maximization

the values of  $\hat{\alpha}$  and  $\hat{\beta}$  are obtained by maximizing the marginal likelihood function which is given by

$$p(\mathbf{t}|\alpha, \beta) = \int p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)d\mathbf{w} \quad (89)$$

where

$$p(\mathbf{t}|\mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (90)$$

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

We can re-estimate  $\alpha$  and  $\beta$  analytically or algorithmically.

- Analytical way: evaluate the evidence function, set its derivative equal to zero and then obtain re-estimation equations for  $\alpha$  and  $\beta$
- Alternatively, use the technique called the expected maximization algorithm.

# The Evidence Approximation

## Evaluate the evidence function analytically

The marginal likelihood is obtained by integrating out parameters

$$p(\mathbf{t}|\alpha, \beta) = \int p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)d\mathbf{w} \quad (91)$$

and then we can write the evidence function, by (90), in the form

$$p(\mathbf{t}|\alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp\{-E(\mathbf{w})\}d\mathbf{w} \quad (92)$$

where

$$\begin{aligned} E(\mathbf{w}) &= \beta E_D(\mathbf{w}) + \alpha E_W(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \Phi\mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \\ &= E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N) \end{aligned} \quad (93)$$

where

$$\mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi \quad \mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t} \quad (94)$$



# The Evidence Approximation

Hence, we get

$$\int \exp \{-E(\mathbf{w})\} d\mathbf{w} = \exp \{-E(\mathbf{m}_N)\} (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \quad (95)$$

Taking logarithm to (91), using (95), the log of marginal likelihood is given by

$$\ln p(\mathbf{t}|\alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) - \frac{1}{2} \ln |\mathbf{A}| - \frac{N}{2} \ln(2\pi) \quad (96)$$



# The Evidence Approximation

## Set its derivative equal to zero (WRT $\alpha$ )

Defining the following eigenvector equation

$$(\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (97)$$

then  $\mathbf{A}$  has eigenvalues  $\alpha + \lambda_i$  by

$$\mathbf{A} \mathbf{u}_i = (\alpha + \lambda_i) \mathbf{u}_i \quad (98)$$

Evaluate the derivative of the log of the marginal likelihood

$$\begin{aligned} \frac{d}{d\alpha} \ln p(\mathbf{t}|\alpha, \beta) &= \frac{d}{d\alpha} \left( \frac{M}{2} \ln \alpha \right) - \frac{d}{d\alpha} (\ln |\mathbf{A}|) - \frac{d}{d\alpha} E(\mathbf{m}_N) \\ &\simeq \frac{M}{2\alpha} - \frac{d}{d\alpha} \frac{1}{2} \ln \mathbf{A} - \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N \end{aligned} \quad (99)$$

Since

$$\frac{d}{d\alpha} \ln |\mathbf{A}| = \frac{d}{d\alpha} \ln \prod_{i=1}^M (\lambda_i + \alpha) = \sum_{i=1}^M \frac{1}{\lambda_i + \alpha} \quad (100)$$

then the stationary point is

$$0 = \frac{M}{2\alpha} - \frac{1}{2} \sum_{i=1}^M \frac{1}{\lambda_i + \alpha} - \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (101)$$



# The Evidence Approximation

Rearrange the equation, then we get

$$\alpha \mathbf{m}_N^T \mathbf{m}_N = M - \alpha \sum_{i=1}^M \frac{1}{\lambda_i + \alpha} = \gamma \quad (102)$$

Since  $M = \sum_{i=1}^M 1$ , then

$$\alpha \mathbf{m}_N^T \mathbf{m}_N = \sum_{i=1}^M \frac{\lambda_i}{\lambda_i + \alpha} = \gamma \quad \text{or} \quad \alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} \quad (103)$$

## Rearrange sequence

- 1 Set prior distribution  $p(\mathbf{w})$  over parameters, get  $\alpha$ , calculate  $\gamma$
- 2 Evaluate posterior distribution  $p(\mathbf{w}|\mathbf{t})$ , get  $\mathbf{m}_N$
- 3 Re-estimate  $\alpha$  as  $\gamma / \mathbf{m}_N^T \mathbf{m}_N$



# The Evidence Approximation

**Set its derivative equal to zero (WRT  $\beta$ )**

$$\frac{d}{d\beta} \ln |\mathbf{A}| = \frac{d}{d\beta} \sum_i \ln(\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta} \quad (104)$$

Then, the stationary point is given by solving the following equation

$$0 = \frac{N}{2\beta} - \frac{1}{2}(\mathbf{t} - \Phi \mathbf{m}_N)^T (\mathbf{t} - \Phi \mathbf{m}_N) - \frac{\gamma}{2\beta} \quad (105)$$

Therefore, solving the equation, we get

$$\beta^{-1} = \frac{1}{N - \gamma} \sum_{i=1}^N \{t_n - \mathbf{m}_N^T \phi(\mathbf{x}_n)\}^2 \quad (106)$$





## An Elegant Interpretation of $\alpha$ and $\gamma$

Since,  $\Phi^T \Phi$  is positive definitive, the eigenvector equation must have positive eigenvalues  $\lambda_i \geq 0$ . And then we get

$$0 \leq \frac{\lambda_i}{\lambda_i + \alpha} \leq 1 \quad (107)$$

And then, we may consider two cases of the value of  $\lambda_i$

■ ( $\lambda_i \gg \alpha$ ):

The corresponding parameter  $w_i$  will be close to its maximum likelihood value.

The value of  $w_i$  is tightly constrained by the data. Such parameters are called *well determined*

■ ( $\lambda_i \ll \alpha$ ):

The corresponding parameter  $w_i$  will be close to zero. The value of  $w_i$  is loosely constrained by the data.

Now, we can construct a feasible range of  $\gamma$  as follows

$$0 \leq \gamma = \sum_{i=1}^M \frac{\lambda_i}{\lambda_i + \alpha} \leq M \quad (108)$$

As the extension of the concept *well determined*, the quantity  $\gamma$  measures the effective total number of well determined parameters.