

## 1.2.6 Bayesian curve fitting

Note: the Gaussian curve problem demonstrated, can be/is solved analytically

In a fully Bayesian approach, we should consistently apply the sum and product rules of probability, which requires, as we shall see shortly, that we integrate over all values of  $w$ . Such marginalizations lie at the heart of Bayesian methods for pattern recognition.

given the training data  $\mathbf{x}$  and  $\mathbf{t}$ , along with a new test point  $x$ , and our goal is to predict the value of  $t$ .

we assume that the parameters  $\alpha$  and  $\beta$  are fixed and known in advance (in later chapters we shall discuss how such parameters can be inferred from data in a Bayesian setting).

predictive  
distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}. \quad (1.68)$$

is the posterior distribution

can be found by normalizing the right-hand side of (1.66).

omitted the dependence on  $\alpha$  and  $\beta$  to simplify the notation

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha). \quad (1.66)$$

posterior  $\propto$  likelihood  $\times$  prior

target values  $\mathbf{t} = (t_1, \dots, t_N)^T$ . (sometimes  $t$  is labels)

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j \quad (1.1)$$

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \quad (1.60)$$

$\beta$  corresponding to the inverse variance of the distribution.

$N$  input values  $\mathbf{x} = (x_1, \dots, x_N)^T$  where  $x$  are observations or predictions

## 1.2.6 Bayesian curve fitting (cont.)

solving 1.60:

■ solve for  $\mathbf{w}_{\text{ML}}$  by minimizing the negative log-likelihood:

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi). \quad (1.62)$$

$$\frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 = 0$$

$$\text{where } y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

One can use regularization to avoid overfitting by a polynomial order that is too high:

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where  $\|\mathbf{w}\|^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$ , and the coefficient  $\lambda$  governs the relative importance of the regularization term compared with the sum-of-squares error term. Note that often the coefficient  $w_0$  is omitted from the regularizer because its

■ solve for  $\beta_{\text{ML}}$

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2. \quad (1.63)$$

## 1.2.6 Bayesian curve fitting (cont.)

then 1.60 becomes 1.64:

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \quad (1.60)$$

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1}). \quad (1.64)$$

for the prior of 1.66 consider a Gaussian over the polynomial coefficients of  $\mathbf{w}$ :

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\} \quad (1.65)$$

where  $\alpha$  is the precision of the distribution, and  $M+1$  is the total number of elements in the vector  $\mathbf{w}$  for an  $M^{\text{th}}$  order polynomial. Variables such as  $\alpha$ , which control

■ solve for  $\mathbf{w}$  by minimizing the negative log-likelihood of 1.66:

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}. \quad (1.67)$$

Thus we see that maximizing the posterior distribution is equivalent to minimizing the regularized sum-of-squares error function encountered earlier in the form (1.4), with a regularization parameter given by  $\lambda = \alpha/\beta$ .

Integration of 1.68 can be solved analytically:

**predictive distribution**

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}. \quad (1.68)$$

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \quad (1.69)$$

where the mean and variance are given by

$$m(x) = \beta \phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(x_n) t_n \quad (1.70)$$

$$s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x). \quad (1.71)$$

Here the matrix  $\mathbf{S}$  is given by

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \quad (1.72)$$

where  $\mathbf{I}$  is the unit matrix, and we have defined the vector  $\phi(x)$  with elements  $\phi_i(x) = x^i$  for  $i = 0, \dots, M$ .

$\mathbf{x}$  = 10 data points drawn randomly from 0 to 1.

$y$  = the values generated from  $\sin(2\pi x)$

or use those from figures:

$x = (0., 6./60, 13./60, 21./60, 27./60, 34./60, 40./60, 47./60, 54./60, 60./60)$

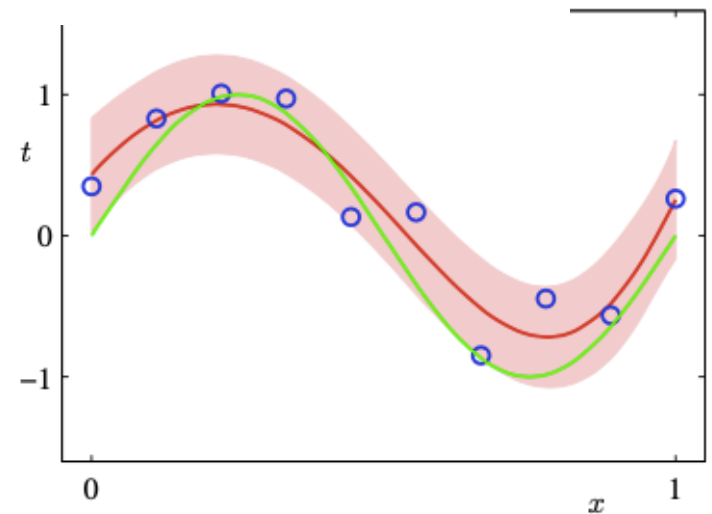
$y = (3.5/9, 7.5/9, 9./9, 8.5/9, -1./9, -1.5/9, -8./9, -4./9, -5./9, 2.5/9)$

alpha = 5E-3

beta = 11.1

**Figure 1 17**

The predictive distribution resulting from a Bayesian treatment of polynomial curve fitting using an  $M = 9$  polynomial, with the fixed parameters  $\alpha = 5 \times 10^{-3}$  and  $\beta = 11.1$  (corresponding to the known noise variance), in which the red curve denotes the mean of the predictive distribution and the red region corresponds to  $\pm 1$  standard deviation around the mean.



where  $\mathbf{I}$  is the unit matrix, and we have defined the vector  $\phi(x)$  with elements  $\phi_i(x) = x^i$  for  $i = 0, \dots, M$ .

$\phi_i(x) = x^i$ . let  $x=x_0=0$ . then  $\phi_i(x_0) = 0^i$ .  $\phi(x_0)=(0^0,0^1,0^2,...0^9)$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(x_n) \phi(x)^T \tag{1.72}$$

$\mathbf{x}$  = 10 data points drawn randomly from 0 to 1.  
 $y$  = the values generated from  $\sin(2\pi x)$   
 or use those from figures:  
 $\mathbf{x} = (0., 6./60, 13./60, 21./60, 27./60, 34./60, 40./60, 47./60, 54./60, 60./60)$   
 $y = (3.5/9, 7.5/9, 9./9, 8.5/9, -1./9, -1.5/9, -8./9, -4./9, -5./9, 2.5/9)$   
 $\alpha = 5E-3$   
 $\beta = 11.1$   
 $M$  = order of fit = 9 for this example

$\phi_i(x) = x^i$ . let  $x=x_0=0$ . then  $\phi_i(x_0) = 0^i$ .  $\phi(x_0)=(0^0,0^1,0^2,...0^9)$   
 1 to N are data element numbers from 1 to 10 (or 0-9)  
 0 to M are the order (power exponents)

$\phi_0(x) =$	$0.^0$		$\phi(x_0)^T=(0^0,0^1,0^2,...0^9)$
$[N \times 1]$	$6./60^0$		$\phi(x_1)^T=((6./60)^0,(6./60)^1,...(6./60)^9)$
	$13./60^0$		$\phi(x_9)^T=((60./60)^0,(60./60)^1,...(60./60)^9)$
	$21./60^0$		
	$27./60^0$		
	$34./60^0$		
	$40./60^0$		
	$47./60^0$		
	$54./60^0$		
	$60./60^0$		
$\phi(x) =$	$0.^0, \dots, 0.^9$		
$[N \times M]$	$6./60^0, \dots, 6./60^9$		
	$13./60^0, \dots, 13./60^9$		
	$21./60^0, \dots, 21./60^9$		
	$27./60^0, \dots, 27./60^9$		
	$34./60^0, \dots, 34./60^9$		
	$40./60^0, \dots, 40./60^9$		
	$47./60^0, \dots, 47./60^9$		
	$54./60^0, \dots, 54./60^9$		
	$60./60^0, \dots, 60./60^9$		

each is  $[1 \times M]$

search literature. Vectors are denoted by lower case bold Roman letters such as  $\mathbf{x}$ , and all vectors are assumed to be column vectors. A superscript  $T$  denotes the transpose of a matrix or vector, so that  $\mathbf{x}^T$  will be a row vector. Uppercase bold roman letters, such as  $\mathbf{M}$ , denote matrices. The notation  $(w_1, \dots, w_M)$  denotes a row vector with  $M$  elements, while the corresponding column vector is written as  $\mathbf{w} = (w_1, \dots, w_M)^T$ .

If we have  $N$  values  $\mathbf{x}_1, \dots, \mathbf{x}_N$  of a  $D$ -dimensional vector  $\mathbf{x} = (x_1, \dots, x_D)^T$ , we can combine the observations into a data matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row of  $\mathbf{X}$  corresponds to the row vector  $\mathbf{x}_n^T$ . Thus the  $n, i$  element of  $\mathbf{X}$  corresponds to the  $i^{\text{th}}$  element of the  $n^{\text{th}}$  observation  $\mathbf{x}_n$ . For the case of one-dimensional variables we shall denote such a matrix by  $\mathbf{x}$ , which is a column vector whose  $n^{\text{th}}$  element is  $x_n$ . Note that  $\mathbf{x}$  (which has dimensionality  $N$ ) uses a different typeface to distinguish it from  $x$  (which has dimensionality  $D$ ).

here,  $\mathbf{x}$  is a row vector

training data comprising  $N$  input values  $\mathbf{x} = (x_1, \dots, x_N)^T$

original data set consists of  $N$  data points  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \overbrace{\phi(x_n)^T \phi(x)}^{\phi(x)^T * \phi(x)} \tag{1.72}$$

$$m(x) = \beta \cancel{\phi(x)^T} \mathbf{S} \sum \phi(x_n) t_n \tag{1.70}$$

~~$[M \times N][M \times M][M \times 1] = [M \times 1]$~~  error in eqn. see ctgk PRML

$$s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x). \tag{1.71}$$

~~$[[M \times M]] = [M \times M]$~~  error in eqn. see ctgk PRML github

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \tag{1.69}$$