

---

# Reports on PRML Reading Talk

---

Siyu Wang  
Department of Computer Science  
Tsinghua University  
thuwangsy@gmail.com

## 1 Introduction to ML, probability basics–prml chap.1,2, mlapp chap.1,2

The key purpose of machine learning is to extract good **features** and **patterns** from **data**. To realize this purpose, we mainly use function fitting as the tool to solve this problem. In this section we will focus on a polynomial fitting problem, solving it from different angles and considering differences and relations among them.

But we must keep in mind that, function fitting is only a mathematical tool, not our purpose.

Consider the following problem: given some observed points:

$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

in which,  $t = \sin(2\pi x) + \mathcal{N}(0, \sigma^2)$ . Then we want to find a polynomial function

$$t = y(x, \mathbf{w}) = \mathbf{w}_0 + \mathbf{w}_1 x + \mathbf{w}_2 x^2 + \dots + \mathbf{w}_M x^M \quad (1.1)$$

to fit these points, so when given new values of  $x$ , we can predict the corresponding  $t$ .

### From function fitting angle

We want to optimize the unknown parameter  $\mathbf{w}$  in equation.1.1 by minimize the *error function*

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2. \quad (1.2)$$

So we can set the error function's derivative with respect to  $\mathbf{w}$  to zero and easily get the optimal parameter  $\mathbf{w}^*$ .

To reduce the over-fitting problem, we may make a little change about the error function by adding a regularization item:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2. \quad (1.3)$$

### From probabilistic angle

Here, we shall assume that, given the value of  $x$ , the corresponding value of  $t$  has a Gaussian distribution with a mean equal to the value  $y(x, \mathbf{w})$  and a fixed variance  $\beta$ . So we have

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(y(x, \mathbf{w}), \beta^{-1}) \quad (1.4)$$

### MLE: maximum likelihood estimation

Given  $N$  input value  $\mathbf{x} = (x_1, \dots, x_N)^T$  and their corresponding target value  $\mathbf{t} = (t_1, \dots, t_N)^T$ , we have the likelihood function with the form as

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}). \quad (1.5)$$

It's convenient to maximize the logarithm of the likelihood function

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi). \quad (1.6)$$

So when determining  $\mathbf{w}$ , maximizing likelihood function is equal to minimizing the error function in equation.1.2. And we obtain

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{ML}) - t_n\}^2. \quad (1.7)$$

So we get a probability distribution of  $t$  when given a new value of  $x$

$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t|y(x, \mathbf{w}_{ML}), \beta_{ML}) \quad (1.8)$$

### MAP: maximum a posteriori estimation

Now we take a step towards a more Bayesian approach and introduce a prior distribution over the polynomial coefficient  $\mathbf{w}$  and for simplicity, consider a Gaussian distribution of the form

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\|\mathbf{w}\|^2\right\} \quad (1.9)$$

where  $\alpha$  is the precision of the distribution and  $M + 1$  is the total number of elements in the vector  $\mathbf{w}$ . Using Bayes' theorem, we can get the posterior distribution for  $\mathbf{w}$ , which is proportional to the product of the prior distribution and the likelihood function

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{w}, \mathbf{x}, \beta)p(\mathbf{w}|\alpha). \quad (1.10)$$

And finally we find that the maximum of the posterior is given by the minimum of

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}. \quad (1.11)$$

Thus we see that maximizing the posterior distribution is equivalent to minimizing the regularized sum-of-square error function in the form (1.3)

### Inference and decision

Here, we come to a classification problem and we have three different approaches to doing inference and decision.

1. First solve  $p(\mathbf{x}|\mathcal{C}_k)$  as well as  $p(\mathcal{C}_k)$  for each class individually, then figure out  $p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$  and  $p(\mathbf{x})$  can be gotten from  $p(\mathbf{x}) = \sum_k p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ . This is equivalent to model the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$
2. Determine posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$  and then use decision theory. *discriminative model*.
3. Find  $f(\mathbf{x})$  which maps  $\mathbf{x}$  directly to a class label.

The first approach is too time-consuming and demanding, while the last one is simplest but not so robust as approach.2 because we can see a lot of information from the posterior probabilities.

## **For others**

Other than the main ideas and problems discussed above, there are still some other small items here: *Gaussian probability distribution, exponential family, students' t distribution, conjugate prior*, and so on. For more details about these, we can refer to the original book.

- 2 Statistics, information theory basics–mlapp chap.5,6**
- 3 Linear models(1)–mlapp chap.7, prml chap.3**
- 4 Linear models(2)–mlapp chap.8, prml chap.4**
- 5 Kernels–mlapp chap.14. prml chap.6-7**
- 6 Graphical models–mlapp chap 10, 19. prml chap 8.1-8.3**
- 7 Latent variable model, clustering and EM. prml chap9. mlapp chap11, 25**
- 8 Continuous latent variables. prml chap 12. mlapp 12-13**
- 9 Exact inference–mlapp chap20. prml chap 8.4**
- 10 Variational inference–prml chap10. mlapp chap21-22**
- 11 Monte Carlo– prml11. mlapp 23,24**
- 12 Sequential data–prml chap13. mlapp chap 17-18.**
- 13 Gaussian Process–mlapp chap15**
- 14 Neural network–mlapp chap 16. prml chap5.**
- 15 Deep learning–mlapp chap 28.**
- 16 Combining models. prml chap 14**