

1 Introduction Problem Setting

The introduction use polynomial curve fitting as example. We can use the general defined loss function to describe the fitting performance and the normalized RMS loss function.

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

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

We can use the penalty term to control the model complexity in order to match the problem complexity. (As mentioned in class, the DL model can be overparameterized model, which does not follow the traditional model complexity control theory)

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

2 Probability Theory

This section describes the general probability definition in discrete conditions and continuous conditions. With Bayes' theorem and the distribution, we can describe the learning procedure from a probability perspective. ¹

3 The curve fitting revisited in probabilistic perspective

Given the value of x , the corresponding value of t has a Gaussian distribution with a mean equal to the value $y(x, \mathbf{w})$

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

, where precision parameter β corresponding to the inverse variance of the distribution.

If the data point is i.i.d.,

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

thus we can have the log likelihood,

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

We can remove the terms that do not depend on \mathbf{w} , and scale the whole equation, then we can minimize the negative log likelihood,

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

For inference, we can have, $p(t | x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t | y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$ and the Prediction Distribution, $p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t | x, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$.

4 Model Selection

Cross-validation drawbacks: the training is itself computationally expensive and has multiple complexity parameters for a single model.

¹Limited to the page limit, you can get the full version note at https://github.com/Xiang-Pan/NYU_Baysian_Machine_Learning/blob/master/reading_notes/ch1/note1.pdf

5 The Curse of Dimensionality

We have some observations and assumptions for the high-dimensional data, which helps us overcome the curse of dimensionality.

6 Decision Theory

Bayes' theorem, $p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$. Minimizing the misclassification rate (or maximizing the correct classification rate) to make the decision. (PS, these two inference methods can be the same at the inference time. However, they MAY be different on the training procedure)

7 Inference and decision

We have three typical ways to describe the inference procedure, however, we may prefer those methods which can keep the posterior probability.