

1 Introduction Problem Setting

The introduction use polynomial curve fitting as the example.

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

we can use the general defined loss function to describe the fitting performance and the normalized RMS loss function(in order to make compare different dataset size).

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

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

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 describe the general probability definition in discrete condition and continuous condition. With Bayes' theorem and the distribution, we can describe the learning process in probability perspective.

1

3 The curve fitting revisited in probabilistic perspective

3.1 Prior

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.

¹limited to the page limit, you can get the full version note at https://github.com/Xiang-Pan/NYU_Bayesian_Machine_Learning

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

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

thus we can have the log likelihood,

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

We can remove the term 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,

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

4 Model Selection

cross-validation drawbacks:

- training is itself computationally expensive
- multiple complexity parameters for a single model

5 The Curse of Dimensionality

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

6 Decision Theory

Bayes' theorem,

$$p(\mathcal{C}_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathcal{C}_k) p(\mathcal{C}_k)}{p(\mathbf{x})}$$

Minimizing the misclassification rate,

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} \end{aligned}$$

We can maximize the probability of being correct

$$\begin{aligned} p(\text{ correct }) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\ &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \end{aligned}$$

(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

Inference stage in which we use training data to learn a model for $p(C_k|x)$.

Decision stage in which we use these posterior probabilities to make optimal class assignments.

Three ways,

- Bayes' theorem
- Directly assign class for each x
- Discriminant function

8 Information Theory

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x} | \mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y} | \mathbf{x}]$$

Thus we can view the mutual information as the reduction in the uncertainty about x by virtue of being told the value of y (or vice versa). From a Bayesian perspective, we can view $p(x)$ as the prior distribution for x and $p(x|y)$ as the posterior distribution after we have observed new data y . The mutual information therefore represents the reduction in uncertainty about x as a consequence of the new observation y .