



National Yang Ming Chiao Tung University

[EECM30060] Machine Learning

Homework 1

Po-Chuan, Chen

Student ID: 311511052

present90308.ee11@nycu.edu.tw

NATIONAL YANG MING CHIAO TUNG UNIVERSITY

October 14, 2022

Contents

1	Bayesian Linear Regression	2
1.1	Question	2
1.2	Proof	3

1 Bayesian Linear Regression

1.1 Question

A linear regression function can be expressed as below where the ϕ is a basis function:

$$y(x, \mathbf{w}) = \mathbf{w}^\top \phi(x)$$

In order to make prediction of t for new test data x from the learned \mathbf{w} , we may

- multiply the likelihood function of new data $p(t|x, \mathbf{w})$ and the posterior distribution of training set with label set.
- take the integral over \mathbf{w} to find the predictive distribution

$$\begin{aligned} p(t|x, \mathbf{x}, \mathbf{t}) &= \int_{-\infty}^{\infty} p(t, \mathbf{w}|x, \mathbf{x}, \mathbf{t}) d\mathbf{w} \\ &= \int_{-\infty}^{\infty} p(t|\mathbf{w}, x, \mathbf{x}, \mathbf{t}) p(\mathbf{w}|x, \mathbf{x}, \mathbf{t}) d\mathbf{w} \\ &= \int_{-\infty}^{\infty} p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}. \end{aligned}$$

Prove that the predictive distribution just mentioned is the same with the form

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

where

$$\begin{aligned} m(x) &= \beta \phi(x)^\top \mathbf{S} \sum_{n=1}^N \phi(x_n) t_n \\ s^2(x) &= \beta^{-1} + \phi(x)^\top \mathbf{S} \phi(x). \end{aligned}$$

1.2 Proof

Based on the marginal and conditional Gaussian, a marginal Gaussian distribution for \mathbf{x} and conditional

Gaussian distribution for \mathbf{y} given \mathbf{x} in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mu, \Lambda^{-1}) \quad (1)$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (2)$$

and the marginal distribution of \mathbf{y} and the conditional distribution of \mathbf{x} given \mathbf{y} are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mu + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T) \quad (3)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(\mathbf{x} \mid \sum \{\mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \Lambda\mu\}, \sum) \quad (4)$$

where

$$\sum = (\Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}$$

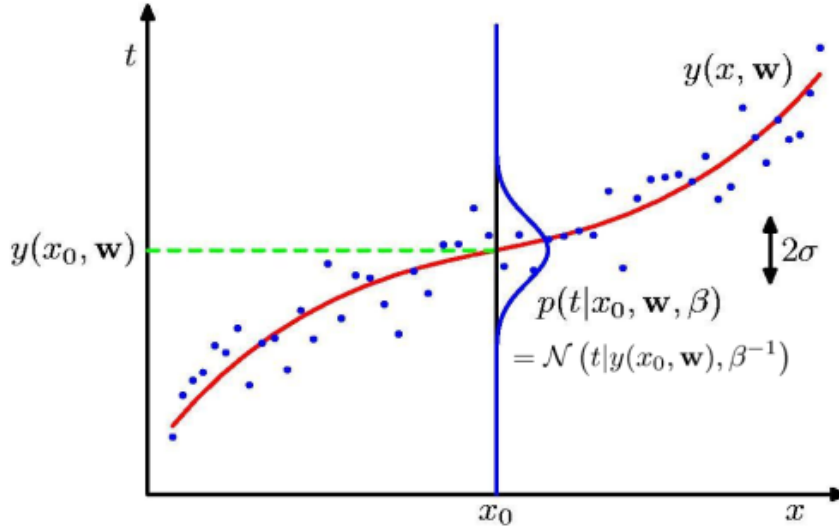


Figure 1: Likelihood model

We set the likelihood as

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

and the posterior as

$$p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} \mid \beta \mathbf{S}_N \sum_{n=1}^N t_n \phi(x_n), \mathbf{S}_N).$$

Based on the problem, we can take the predictive distribution to

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T) \implies$$

$$p(t | x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | \boldsymbol{\phi}(x)^T \beta \mathbf{S}_N \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n), \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x))$$

which

$$\mathbf{x} \leftarrow \mathbf{w}, \mu = \beta \mathbf{S}_N \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n), \Lambda^{-1} = \mathbf{S}_N$$

$$\mathbf{y} \leftarrow t, \mathbf{A} = \boldsymbol{\phi}(x)^T, \mathbf{b} = 0, \mathbf{L}^{-1} = \beta^{-1}$$

and \mathbf{S}_N is covariance matrix with N dimension, α which is precision of prior, and β is precision in the data.

If we want to calculate the predictive distribution in full Bayesian treatment with training data \mathbf{x} and \mathbf{t} and a new test point x , the distribution will turn into

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

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

where the mean $m(x)$ and the variance $s^2(x)$ is like

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

$$s^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$$

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

$$\boldsymbol{\phi}(x_n) = \begin{Bmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{Bmatrix} \text{ with polynomial regression, } \mathbf{I} = \text{unit matrix } M \times M$$

and then we can use predictive distribution to predict t based on the new coming \mathbf{x} .

And the distribution can be transfer into

$$p(t \mid x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t \mid m(x), \sigma_N^2(x))$$

where the mean and variance is

$$m(x) = \boldsymbol{\phi}(x)^T m_N, \quad m_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$$

$$\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \xrightarrow{N \rightarrow \infty} \beta^{-1} \text{ and } \sigma_{N+1}^2(x) \leq \sigma_N^2(x) \quad (5)$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

As we can see in the equation (5), we can know that

$$\mathbf{S}_{N+1}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^{N+1} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T = \mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_N) \boldsymbol{\phi}(x_N)^T$$

, note that we can see $\alpha \rightarrow 0$ implies $m_N \rightarrow w = (\boldsymbol{\phi}^T \boldsymbol{\phi})^{-1} \boldsymbol{\phi}^T \mathbf{t}$, and

$$(\mathbf{M} + \mathbf{v} \mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1} \mathbf{v})(\mathbf{v}^T \mathbf{M}^{-1})}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}}$$

So, based on the results we get, we can know that

$$\begin{aligned} \sigma_{N+1}^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_{N+1} \boldsymbol{\phi}(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T (\mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_N) \boldsymbol{\phi}(x_N)^T) \\ &= \beta^{-1} + \boldsymbol{\phi}(x)^T [\mathbf{S}_N - \frac{\beta (\mathbf{S}_N \boldsymbol{\phi}(x_N)) (\boldsymbol{\phi}(x_N)^T \mathbf{S}_N)}{1 + \beta \boldsymbol{\phi}(x_N)^T \mathbf{S}_N \boldsymbol{\phi}(x_N)} \boldsymbol{\phi}(x_N)] \boldsymbol{\phi}(x) \\ &= \sigma_N^2(x) - \frac{\beta (\boldsymbol{\phi}(x_N)^T \mathbf{S}_N \boldsymbol{\phi}(x))^2}{1 + \beta \boldsymbol{\phi}(x_N)^T \mathbf{S}_N \boldsymbol{\phi}(x_N)} \leq \sigma_N^2(x) \end{aligned}$$

So the predictive distribution used as follows

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

$$p(t \mid x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t \mid m(x), \sigma^2(x))$$

where the mean $m(x)$ and the variance $\sigma^2(x)$ is like

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

$$\sigma^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$$

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

$$\boldsymbol{\phi}(x_n) = \begin{Bmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \phi_2(x_n) \\ \vdots \\ \phi_{M-1}(x_n) \end{Bmatrix}, \quad \mathbf{I} = \text{unit matrix } M \times M,$$

$$\boldsymbol{\phi}(x_n) = \begin{Bmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{Bmatrix} \quad \text{with polynomial regression}$$

which β^{-1} is noise in data, $\boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$ is an uncertainty in \mathbf{w} .