

Homework 1 (Lecture: Bayesian Learning)

Solve Exercises 1–8, and submit an answer sheet at the end of the lecture on 7.1.2016.

1 Linear regression model

Consider a linear regression model with unknown parameter $\mathbf{w} = \mathbf{a} \in \mathbb{R}^M$:

$$p(y|\mathbf{x}, \mathbf{a}) = \text{Norm}_1(y; \mathbf{a}^\top \mathbf{x}, \sigma^2) = \frac{\exp\left(-\frac{(y - \mathbf{a}^\top \mathbf{x})^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}. \quad (1)$$

$\text{Norm}_M(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the M -dimensional Gaussian distribution with mean $\boldsymbol{\mu}$ and the covariance $\boldsymbol{\Sigma}$. We treat the noise variance σ^2 as a fixed constant.

Assume that we observed N samples $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$, and that, for each given input $\mathbf{x}^{(n)}$, the output $y^{(n)}$ was independently and identically (i.i.d.) drawn from $\text{Norm}_1(y; \mathbf{a}^{*\top} \mathbf{x}, \sigma^2)$ with unknown \mathbf{a}^* .

Define

$$\mathbf{y} = (y^{(1)}, \dots, y^{(N)})^\top \in \mathbb{R}^N, \quad \mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)})^\top \in \mathbb{R}^{N \times M}.$$

Then, the model likelihood is written as

$$p(\mathbf{y}|\mathbf{X}, \mathbf{a}) = \text{Norm}_N(\mathbf{y}; \mathbf{X}\mathbf{a}, \sigma^2 \mathbf{I}_N) = \frac{\exp\left(-\frac{\|\mathbf{y} - \mathbf{X}\mathbf{a}\|^2}{2\sigma^2}\right)}{(2\pi\sigma^2)^{N/2}}, \quad (2)$$

where \mathbf{I}_N is the $N \times N$ identity matrix. We adopt Gaussian prior with mean $\mathbf{0}$ and covariance \mathbf{C} .

$$p(\mathbf{a}|\mathbf{C}) = \text{Norm}_M(\mathbf{a}; \mathbf{0}, \mathbf{C}) = \frac{\exp\left(-\frac{1}{2}\mathbf{a}^\top \mathbf{C}^{-1} \mathbf{a}\right)}{(2\pi)^{M/2} |\mathbf{C}|^{1/2}}. \quad (3)$$

2 Posterior and Predictive distributions

Exercise 1: Derive the posterior distribution $p(\mathbf{a}|\mathbf{y}, \mathbf{X}, \mathbf{C})$ on \mathbf{a} .

Exercise 2: Derive the predictive distribution $p(y^*|\mathbf{x}^*, \mathbf{y}, \mathbf{X}, \mathbf{C})$ on the output y^* for a new input \mathbf{x}^* .

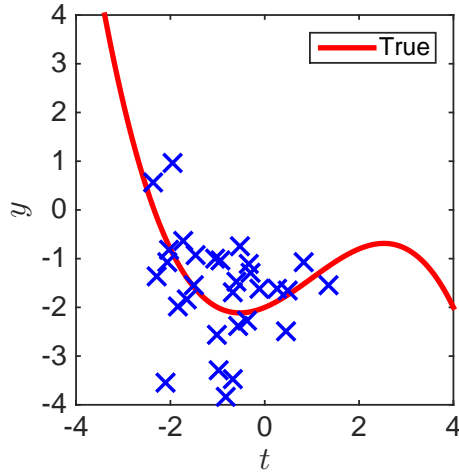


Figure 1: $N = 30$ samples from the linear regression model $y = \mathbf{a}^{*\top} \mathbf{x} + \varepsilon$, where $\mathbf{a}^* = (-2, 0.4, 0.3, -0.1)^\top$, $\mathbf{x} = (1, t, t^2, t^3)^\top$, and $\varepsilon \sim \text{Norm}_1(0, 1^2)$.

Exercise 3: A set of training samples (shown as crosses in Fig.1) are given in a separate file “data.txt”. The upper row corresponds to t and the lower row corresponds to y . Draw the same figure as Fig. 1. Numerically compute the mean \hat{y} and the covariance $\hat{\sigma}_y^2$ of the predictive distribution $p(y^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X}, \mathbf{C})$ as a function of t^* (on a grid, for example, $t^* = -4.00, -3.99, -3.98, \dots, 4.00$). Here, set $\mathbf{C} = 10000 \mathbf{I}_M$ and $\sigma^2 = 1$. Then, overlap the three curves corresponding to \hat{y} and $\hat{y} \pm \hat{\sigma}_y$ in the previous figure. Any programming language or software can be used. Only submit a printed figure.

3 Marginal Likelihood and Empirical Bayesian Learning

Exercise 4: Compute the marginal likelihood $p(\mathcal{D} | \mathbf{C})$ (without omitting any constant factor).

Estimating the hyperparameter \mathbf{C} by maximizing the marginal likelihood is called *empirical Bayesian learning*. Equivalently, we minimize the negative log of the marginal likelihood

$$F^* = -\log p(\mathcal{D} | \mathbf{C}), \quad (4)$$

which is called the *Bayes free energy* or *stochastic complexity*. In addition, $\log p(\mathcal{D} | \mathbf{C})$ is called the *log marginal likelihood* or *evidence*.

Assume that

$$\mathbf{C} = \mathbf{Diag}(c_1^2, \dots, c_M^2) \in \mathbb{D}_{++}^M, \quad \mathbf{X} = \mathbf{I}_M, \quad (5)$$

where \mathbb{D}_{++}^M denotes the set of positive definite diagonal matrices (i.e., $c_m^2 > 0, \forall m$). With the diagonal covariance to be estimated via empirical Bayesian learning, the prior

$$p(\mathbf{a}|\mathbf{C}) = \prod_{m=1}^M \frac{1}{\sqrt{2\pi c_m^2}} \exp\left(-\frac{a_m^2}{2c_m^2}\right)$$

is called *automatic relevance determination* (ARD) prior.

Under the assumption (5), the Bayes free energy can be decomposed as

$$2F^* = \sum_{m=1}^M 2F_m^* + \text{const.}, \quad (6)$$

where each F_m^* depends on c_m^2 but not on $c_{m'}^2$, for $m' \neq m$.

Note that the latter assumption in (5) is only for simplifying the subsequent computation. In typical applications, \mathbf{X} is not diagonal.

Exercise 5: Compute F_m^* .

Exercise 6: Draw F_m^* as a function of c_m^2 for $y_m^2 = 0, 1, 1.5, 2$ and $\sigma^2 = 1$. (Submit a printed figure.)

Exercise 7: Prove that the solution to $\min_{\mathbf{C}} F^*$ is given by

$$\hat{c}_m^2 = \begin{cases} y_m^2 - \sigma^2 & \text{if } y_m^2 > \sigma^2, \\ +0 & \text{otherwise.} \end{cases} \quad (7)$$

Exercise 8: Derive the empirical Bayesian estimator $\hat{\mathbf{a}}^{\text{EB}}$ (the posterior mean of \mathbf{a} with the hyperparameter \mathbf{C} replaced with its estimator).