

ANALYTIC APPROXIMATION OF MARGINAL LIKELIHOOD

18

CHAPTER CONTENTS

Laplace Approximation	197
Approximation with Gaussian Density	197
Illustration	199
Application to Marginal Likelihood Approximation	200
Bayesian Information Criterion (BIC)	200
Variational Approximation	202
Variational Bayesian EM (VBEM) Algorithm	202
Relation to Ordinary EM Algorithm	203

As discussed in Section 17.2, the use of conjugate priors allows us to avoid the explicit computation of the integration in the marginal likelihood:

$$\text{ML}(\beta) = \int \prod_{i=1}^n q(\mathbf{x}_i | \theta) p(\theta; \beta) d\theta.$$

In this chapter, general choices of prior probabilities are considered and analytic approximation methods of the marginal likelihood are introduced.

18.1 LAPLACE APPROXIMATION

In this section, *Laplace approximation* is introduced, which allows us to analytically approximate the integration of any twice-differentiable non-negative function $f(\mathbf{x})$:

$$\int f(\theta) d\theta.$$

18.1.1 APPROXIMATION WITH GAUSSIAN DENSITY

Let $\hat{\theta}$ be the maximizer of $f(\theta)$ with respect to θ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} f(\theta).$$

Let us apply the *Taylor series expansion* (Fig. 2.8) to $\log f(\boldsymbol{\theta})$ about the maximizer $\hat{\boldsymbol{\theta}}$:

$$\begin{aligned} \log f(\boldsymbol{\theta}) &= \log f(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) + \dots, \end{aligned} \quad (18.1)$$

where \mathbf{H} is the *Hessian matrix*:

$$\mathbf{H} = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log f(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}.$$

Since $\hat{\boldsymbol{\theta}}$ is the maximizer of $\log f(\boldsymbol{\theta})$,

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log f(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = \mathbf{0}_b$$

holds and thus the first-order term in Eq. (18.1) is zero.

Let $\log \hat{f}(\boldsymbol{\theta})$ be Eq. (18.1) up to the second-order terms (Fig. 18.1):

$$\log \hat{f}(\boldsymbol{\theta}) = \log f(\hat{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}).$$

Exponentiating both sides yields

$$\hat{f}(\boldsymbol{\theta}) = f(\hat{\boldsymbol{\theta}}) \exp \left(\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right).$$

Recalling that the integration of the normal density is 1,

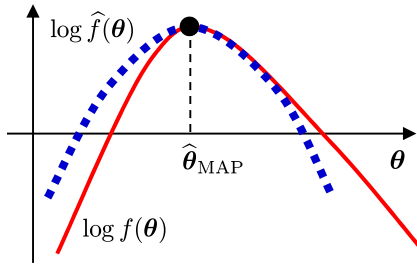
$$\frac{1}{(2\pi)^{\frac{b}{2}} \det(-\mathbf{H})^{\frac{1}{2}}} \int \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top (-\mathbf{H}) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right) d\boldsymbol{\theta} = 1,$$

integration of $\hat{f}(\boldsymbol{\theta})$ yields

$$\int \hat{f}(\boldsymbol{\theta}) d\boldsymbol{\theta} = f(\hat{\boldsymbol{\theta}}) \sqrt{\frac{(2\pi)^b}{\det(-\mathbf{H})}} \approx \int f(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

where b denotes the dimensionality of $\boldsymbol{\theta}$. This is Laplace approximation to $\int f(\boldsymbol{\theta}) d\boldsymbol{\theta}$.

Since approximating $\log f(\boldsymbol{\theta})$ by a quadratic function corresponds to approximating $f(\boldsymbol{\theta})$ by an unnormalized Gaussian function, Laplace approximation is quite accurate if $f(\boldsymbol{\theta})$ is close to Gaussian. For this reason, Laplace approximation is also referred to as *Gaussian approximation*.

**FIGURE 18.1**

Laplace approximation.

18.1.2 ILLUSTRATION

Let us illustrate how to compute Laplace approximation to $\int f(\theta)d\theta$, where

$$f(\theta) = N(\theta; 0, 1^2) + N(\theta; 0, 2^2),$$

$$N(\theta; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\theta - \mu)^2}{2\sigma^2}\right).$$

Note that $N(\theta; \mu, \sigma^2)$ is a normal density and thus the true value is 2.

The maximizer of $f(\theta)$ is given by

$$\hat{\theta} = \operatorname{argmax}_{\theta} f(\theta) = 0.$$

Then

$$\begin{aligned} f(\hat{\theta}) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\hat{\theta}^2}{2}\right) + \frac{1}{\sqrt{8\pi}} \exp\left(-\frac{\hat{\theta}^2}{8}\right) = \frac{3}{2\sqrt{2\pi}}, \\ f'(\hat{\theta}) &= \frac{\partial}{\partial \theta} \log f(\theta) \Big|_{\theta=\hat{\theta}} = -\hat{\theta}N(\hat{\theta}; 0, 1^2) - \frac{\hat{\theta}}{4}N(\hat{\theta}; 0, 2^2) = 0, \\ f''(\hat{\theta}) &= \frac{\partial^2}{\partial \theta^2} \log f(\theta) \Big|_{\theta=\hat{\theta}} = (\hat{\theta}^2 - 1)N(\hat{\theta}; 0, 1^2) + \left(\frac{\hat{\theta}^2}{16} - \frac{1}{4}\right)N(\hat{\theta}; 0, 2^2) \\ &= -\frac{9}{8\sqrt{2\pi}}, \end{aligned}$$

which yield

$$H = \frac{\partial^2}{\partial \theta^2} \log f(\theta) \Big|_{\theta=\hat{\theta}} = \frac{f''(\hat{\theta})f(\hat{\theta}) - f'(\hat{\theta})^2}{f(\hat{\theta})^2} = -\frac{3}{4}.$$

Thus, the Laplace approximation is given by

$$f(0) \sqrt{\frac{2\pi}{-H}} = \sqrt{3} \approx 1.732.$$

18.1.3 APPLICATION TO MARGINAL LIKELIHOOD APPROXIMATION

For the marginal likelihood,

$$\text{ML}(\beta) = \int \prod_{i=1}^n q(\mathbf{x}_i|\theta) p(\theta; \beta) d\theta.$$

Laplace approximation with respect to θ yields

$$\text{ML}(\beta) \approx \prod_{i=1}^n q(\mathbf{x}_i|\hat{\theta}_{\text{MAP}}) p(\hat{\theta}_{\text{MAP}}; \beta) \sqrt{\frac{(2\pi)^b}{\det(-H)}},$$

where b denotes the dimensionality of θ (i.e., the number of parameters), and

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} \left[\sum_{i=1}^n \log q(\mathbf{x}_i|\theta) + \log p(\theta; \beta) \right], \\ H &= \frac{\partial^2}{\partial \theta \partial \theta^\top} \left(\sum_{i=1}^n \log q(\mathbf{x}_i|\theta) + \log p(\theta; \beta) \right) \Big|_{\theta=\hat{\theta}_{\text{MAP}}}. \end{aligned} \quad (18.2)$$

If the number of training samples, n , is large, the *central limit theorem* (see Section 7.4) asserts that the posterior probability $p(\theta|\mathcal{D})$ converges in distribution to the Gaussian distribution. Therefore, the Laplace-approximated marginal likelihood would be accurate when a large number of training samples are available.

18.1.4 BAYESIAN INFORMATION CRITERION (BIC)

The logarithm of the Laplace-approximated marginal likelihood is given by

$$\begin{aligned} \log \text{ML}(\beta) &\approx \sum_{i=1}^n \log q(\mathbf{x}_i|\hat{\theta}_{\text{MAP}}) + \log p(\hat{\theta}_{\text{MAP}}; \beta) \\ &\quad + \frac{b}{2} \log(2\pi) - \frac{1}{2} \log(\det(-H)). \end{aligned} \quad (18.3)$$

Let us further approximate this under the assumption that the number of training samples, n , is large.

The first term $\sum_{i=1}^n \log q(\mathbf{x}_i|\hat{\theta}_{\text{MAP}})$ in Eq. (18.3) has asymptotic order n (see Fig. 14.4):

$$\sum_{i=1}^n \log q(\mathbf{x}_i|\hat{\theta}_{\text{MAP}}) = O(n).$$

On the other hand, the second term $\log p(\hat{\boldsymbol{\theta}}_{\text{MAP}}; \boldsymbol{\beta})$ and the third term $\frac{b}{2} \log(2\pi)$ are independent of n :

$$\log p(\hat{\boldsymbol{\theta}}_{\text{MAP}}; \boldsymbol{\beta}) = O(1) \quad \text{and} \quad \frac{b}{2} \log(2\pi) = O(1).$$

According to the law of large numbers (see Section 7.3), the Hessian matrix \mathbf{H} defined by Eq. (18.2) divided by n converges in probability to $\tilde{\mathbf{H}}$:

$$\frac{1}{n} \mathbf{H} = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \left(\frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i | \boldsymbol{\theta}) + \frac{1}{n} \log p(\boldsymbol{\theta}; \boldsymbol{\beta}) \right) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{MAP}}} \xrightarrow{p} \tilde{\mathbf{H}},$$

where

$$\tilde{\mathbf{H}} = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (E [\log q(\mathbf{x} | \boldsymbol{\theta})]) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{MAP}}}.$$

Since $\tilde{\mathbf{H}}$ is a $b \times b$ matrix,

$$\det(n\tilde{\mathbf{H}}) = n^b \det(\tilde{\mathbf{H}})$$

holds and therefore

$$\frac{1}{2} \log (\det(-\mathbf{H})) \xrightarrow{p} \frac{b}{2} \log n + \frac{1}{2} \log (\det(-\tilde{\mathbf{H}}))$$

is obtained. The first term $\frac{b}{2} \log n$ is proportional to $\log n$, while the second term $\frac{1}{2} \log (\det(-\tilde{\mathbf{H}}))$ is independent of n :

$$\frac{b}{2} \log n = O(\log n) \quad \text{and} \quad \frac{1}{2} \log (\det(-\tilde{\mathbf{H}})) = O(1).$$

Here, suppose that n is large enough so that terms with $O(1)$ can be ignored. Then only $\sum_{i=1}^n \log q(\mathbf{x}_i | \hat{\boldsymbol{\theta}}_{\text{MAP}})$ and $\frac{b}{2} \log n$ remain. Furthermore, the difference between the MAP solution $\hat{\boldsymbol{\theta}}_{\text{MAP}}$ and the maximum likelihood solution $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ is known to be $O(n^{-1})$:

$$\sum_{i=1}^n \log q(\mathbf{x}_i | \hat{\boldsymbol{\theta}}_{\text{MAP}}) = \sum_{i=1}^n \log q(\mathbf{x}_i | \hat{\boldsymbol{\theta}}_{\text{MLE}}) + O(n^{-1}).$$

Consequently, Laplace approximation of the log marginal likelihood given by Eq. (18.3) can be further approximated as

$$\log \text{ML}(\boldsymbol{\beta}) \approx \sum_{i=1}^n \log q(\mathbf{x}_i | \hat{\boldsymbol{\theta}}_{\text{MLE}}) - \frac{b}{2} \log n.$$

The negative of the right-hand side is called the BIC:

$$\text{BIC} = - \sum_{i=1}^n \log q(\mathbf{x}_i | \hat{\boldsymbol{\theta}}_{\text{MLE}}) + \frac{b}{2} \log n.$$

BIC is quite simple and thus is popularly used in model selection. However, BIC is no longer dependent on prior probabilities and thus cannot be used for setting the prior probability. It is also known that BIC is equivalent to the *minimum description length* (MDL) criterion [84], which was derived in a completely different framework.

BIC is similar to AIC explained in Section 14.3, but the second term is different:

$$\text{AIC} = - \sum_{i=1}^n \log q(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_{\text{MLE}}) + b.$$

When $n > e^2 \approx 7.39$, BIC has a stronger penalty than AIC and thus a simpler model would be chosen. However, since AIC and BIC are derived in completely different frameworks (KL divergence approximation and the marginal likelihood approximation), it cannot be simply concluded which one is more superior than the other.

18.2 VARIATIONAL APPROXIMATION

When the integrand $\prod_{i=1}^n q(\mathbf{x}_i | \boldsymbol{\theta}) p(\boldsymbol{\theta}; \boldsymbol{\beta})$ is not close to Gaussian, Laplace-approximated marginal likelihood may not be accurate. In this section, *variational approximation* is introduced, which finds the best approximation to the marginal likelihood in a limited function class that is easier to compute [71].

18.2.1 VARIATIONAL BAYESIAN EM (VBEM) ALGORITHM

The marginal likelihood can be expressed as

$$\text{ML}(\boldsymbol{\beta}) = p(\mathcal{D}; \boldsymbol{\beta}) = \iint p(\mathcal{D}, \boldsymbol{\eta}, \boldsymbol{\theta}; \boldsymbol{\beta}) d\boldsymbol{\eta} d\boldsymbol{\theta},$$

where $\boldsymbol{\eta}$ is called a *latent variable*. Let us consider probability density functions for $\boldsymbol{\eta}$ and $\boldsymbol{\theta}$, denoted by $q(\boldsymbol{\eta})$ and $r(\boldsymbol{\theta})$, called the *trial distributions*. Then Jensen's inequality (see Section 8.3.1) gives the following lower bound:

$$\begin{aligned} \log \text{ML}(\boldsymbol{\beta}) &= \log \iint p(\mathcal{D}, \boldsymbol{\eta}, \boldsymbol{\theta}; \boldsymbol{\beta}) d\boldsymbol{\eta} d\boldsymbol{\theta} \\ &= \log \iint q(\boldsymbol{\eta}) r(\boldsymbol{\theta}) \frac{p(\mathcal{D}, \boldsymbol{\eta}, \boldsymbol{\theta}; \boldsymbol{\beta})}{q(\boldsymbol{\eta}) r(\boldsymbol{\theta})} d\boldsymbol{\eta} d\boldsymbol{\theta} \geq -F(q, r), \end{aligned}$$

where

$$F(q, r) = \iint q(\boldsymbol{\eta}) r(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\eta}) r(\boldsymbol{\theta})}{p(\mathcal{D}, \boldsymbol{\eta}, \boldsymbol{\theta}; \boldsymbol{\beta})} d\boldsymbol{\eta} d\boldsymbol{\theta}$$

is a functional of probability density functions q and r called the *variational free energy* (while $-\log \text{ML}(\beta)$ is called the *free energy*). If q and r are chosen to minimize the variational free energy, a good approximator to the log marginal likelihood may be obtained.

Setting the partial derivatives of the variational free energy at zero,

$$\frac{\partial}{\partial q} F(q, r) = 0 \quad \text{and} \quad \frac{\partial}{\partial r} F(q, r) = 0,$$

yields that the solutions should satisfy

$$q(\eta) \propto \exp \left(\int r(\theta) \log p(\mathcal{D}, \eta | \theta; \beta) d\theta \right), \quad (18.4)$$

$$r(\theta) \propto p(\theta) \exp \left(\int q(\eta) \log p(\mathcal{D}, \eta | \theta; \beta) d\eta \right). \quad (18.5)$$

Since no method is known to analytically solve Eq. (18.4) and Eq. (18.5), these equations are used as updating formulas like the EM algorithm (see Section 15.4). Such an EM-like algorithm is called the VBEM algorithm, and Eq. (18.4) and Eq. (18.5) are, respectively, called the VB-E step and VB-M step.

The variational free energy can be expressed by using the KL divergence (see Section 14.2) as

$$F(q, r) = \text{KL}(q(\eta)r(\theta) \| p(\eta, \theta | \mathcal{D}; \beta)) - \log \text{ML}(\beta).$$

This shows that reducing the variational free energy $F(q, r)$ corresponds to reducing $\text{KL}(q(\eta)r(\theta) \| p(\eta, \theta | \mathcal{D}; \beta))$, and therefore the VBEM algorithm can be regarded as approximating $p(\eta, \theta | \mathcal{D}; \beta)$ by $q(\eta)r(\theta)$.

For this reason, $r(\theta)$ obtained by the VBEM algorithm may be a good approximation to the posterior probability $p(\theta | \mathcal{D}; \beta)$.

18.2.2 RELATION TO ORDINARY EM ALGORITHM

As explained in Section 15.4, the ordinary EM algorithm maximized a lower bound of the likelihood. Here, it is shown that the ordinary EM algorithm can also be interpreted as a variational approximation. *Jensen's inequality* (Section 8.3.1) yields the following lower bound of the log-likelihood:

$$\begin{aligned} \log p(\mathcal{D} | \theta) &= \log \int p(\mathcal{D}, \eta | \theta) d\eta \\ &= \log \int q(\eta | \mathcal{D}, \theta') \frac{p(\mathcal{D}, \eta | \theta)}{q(\eta | \mathcal{D}, \theta')} d\eta \geq b(q, \theta), \end{aligned}$$

where

$$b(q, \theta) = \int q(\eta | \mathcal{D}, \theta') \log \frac{p(\mathcal{D}, \eta | \theta)}{q(\eta | \mathcal{D}, \theta')} d\eta.$$

$\frac{\partial}{\partial q} b(q, \theta) = 0$ yields $q(\eta|\mathcal{D}, \theta') = p(\eta|\mathcal{D}, \theta)$, which is the E-step of the ordinary EM algorithm. Then finding θ' that satisfies $\frac{\partial}{\partial \theta} b(q, \theta)|_{\theta=\theta'} = \mathbf{0}$ is the M-step of the ordinary EM algorithm.

The relation between the ordinary EM algorithm and the VBEM algorithm can be further elucidated by the use of *Dirac's delta function* $\delta(\cdot)$, which satisfies for any function $g : \mathbb{R} \rightarrow \mathbb{R}$ and any real number κ ,

$$\int_{-\infty}^{\infty} g(\tau) \delta(\kappa - \tau) d\tau = g(\kappa).$$

Thus, *convolution* with Dirac's delta function allows us to extract the value of an arbitrary function $g(\cdot)$ at an arbitrary point κ . Dirac's delta function $\delta(\cdot)$ can be expressed as the limit of the normal density:

$$\delta(\tau) = \lim_{\sigma \rightarrow 0} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\tau^2}{2\sigma^2}\right).$$

For multidimensional $\boldsymbol{\tau} = (\tau_1, \dots, \tau_d)^\top$, Dirac's delta function is defined in an elementwise manner as

$$\delta(\boldsymbol{\tau}) = \delta(\tau_1) \times \dots \times \delta(\tau_d).$$

Dirac's delta function, setting

$$r(\theta') = \delta(\theta' - \theta)$$

in the VB-E step yields

$$q(\eta) \propto p(\mathcal{D}, \eta|\theta) \propto p(\eta|\mathcal{D}, \theta),$$

which agrees with the ordinary E-step.