

Bayesian Asymptotics

Aleksandar Dogandžić

March 28, 2017

Contents

Posterior Approximation Around the MAP Estimate 1

Example: Binomial success probability 2

Asymptotic Normality and Consistency for Bayesian Models 4

READING: [Gelman et al. 2014, §4.1 and App. B].

Posterior Approximation Around the MAP Estimate

EXPAND the posterior distribution $f_{\Theta|X}(\theta | x)$ in Taylor series around the maximum *a posteriori* (MAP) estimate

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{MAP}}(x)$$

and keep the first three terms:

$$\ln f_{\Theta|X}(\theta | x) \approx \ln f_{\Theta|X}(\hat{\theta}_{\text{MAP}} | x) + 0.5(\theta - \hat{\theta}_{\text{MAP}})^T \frac{\partial^2 \ln f_{\Theta|X}(\theta | x)}{\partial \theta \partial \theta^T} \bigg|_{\theta = \hat{\theta}_{\text{MAP}}} (\theta - \hat{\theta}_{\text{MAP}}) \quad (1)$$

The second term in the Taylor-series expansion vanishes because the log posterior probability density function (pdf) (or probability mass function (pmf)) has zero derivative at the MAP estimate, see (11) in handout `multivarBayes`.

If the number of measurements is large¹, the posterior distribution $f_{\Theta|X}(\theta | x)$ will be *unimodal*. Furthermore, if $\hat{\theta}_{\text{MAP}}$ is in the interior of the parameter space sp_{Θ} (preferably *far from the boundary* of sp_{Θ}), then we can use the following approximation for the posterior pdf of θ :

$$f_{\Theta|X}(\theta | x) \approx \mathcal{N}(\theta | \hat{\theta}_{\text{MAP}}, \mathcal{J}^{-1}(\hat{\theta}_{\text{MAP}})).$$

where

$$\mathcal{J}(\theta) \triangleq - \frac{\partial^2 \ln f_{\Theta|X}(\theta | x)}{\partial \theta \partial \theta^T}$$

is the *observed information* [Gelman et al. 2014, §4.1].

This approximation follows by looking at (1) as a function of θ .

☞ This Gaussian approximation is usually good if the posterior distribution is unimodal and roughly symmetric.

¹ and asymptotics works, e.g., we have independent, identically distributed (i.i.d.) observations with usual regularity conditions

We may obtain another (simpler, but typically poorer) approximation by replacing the covariance matrix of the above Gaussian distribution with Cramér-Rao bound (CRB) evaluated at $\hat{\theta}_{\text{MAP}}$. If $f_{\Theta|X}(\theta | x)$ has multiple modes (and we can find them all), it can be approximated by a Gaussian mixture or, more generally, a t -distribution mixture.

Example: Binomial success probability

Assume conjugate prior $\theta \sim \text{Beta}(\alpha, \beta)$. Then the posterior is

$$f_{\Theta|X}(\theta | x) \propto \theta^{\alpha+x-1}(1-\theta)^{\beta+N-x-1}$$

see `handout introBayes`.

The derivatives are

$$\begin{aligned} \frac{\partial \ln f_{\Theta|X}(\theta | x)}{\partial \theta} &= \frac{\alpha + x - 1}{\theta} - \frac{\beta + N - x - 1}{1 - \theta} \\ -\frac{\partial^2 \ln f_{\Theta|X}(\theta | x)}{\partial \theta^2} &= \frac{\alpha + x - 1}{\theta^2} + \frac{\beta + N - x - 1}{(1 - \theta)^2} \end{aligned}$$

which gives

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \frac{\alpha + x - 1}{\alpha + \beta + N - 2} \\ -\frac{\partial^2 \ln f_{\Theta|X}(\hat{\theta}_{\text{MAP}} | x)}{\partial \theta^2} &= -\frac{\partial^2 \ln f_{\Theta|X}(\theta | x)}{\partial \theta^2} \Big|_{\theta=\hat{\theta}_{\text{MAP}}} \\ &= \frac{\alpha + \beta + N - 2}{\hat{\theta}_{\text{MAP}}(1 - \hat{\theta}_{\text{MAP}})} \end{aligned}$$

and

$$f(\theta | x) \approx \mathcal{N}\left(\theta \mid \hat{\theta}_{\text{MAP}}, \frac{\hat{\theta}_{\text{MAP}}(1 - \hat{\theta}_{\text{MAP}})}{\alpha + \beta + N - 2}\right).$$

For $x = 37$ and $N = 41$ with Beta(3, 3) prior, we have

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \frac{3 + 37 - 1}{3 + 3 + 41 - 2} \\ &= \frac{39}{45} = 0.8667 \\ J(\hat{\theta}_{\text{MAP}}) &= \frac{45^3}{39 \cdot 6} = 389.42, \quad [J(\hat{\theta}_{\text{MAP}})]^{-1/2} = 0.0507 \end{aligned}$$

see Fig. 1.

Now with the $\frac{1}{2}$ Beta(8, 2) + $\frac{1}{2}$ Beta(2, 8) mixture prior,

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= 0.8989 \\ J(\hat{\theta}_{\text{MAP}}) &= \frac{45^3}{39 \cdot 6} = 490.11, \quad [J(\hat{\theta}_{\text{MAP}})]^{-1/2} = 0.0452 \end{aligned}$$

see Fig. 2.

The comment about unimodal and symmetric is important. However, when the number of observations gets large, this is usually not a problem. See Fig. 3.

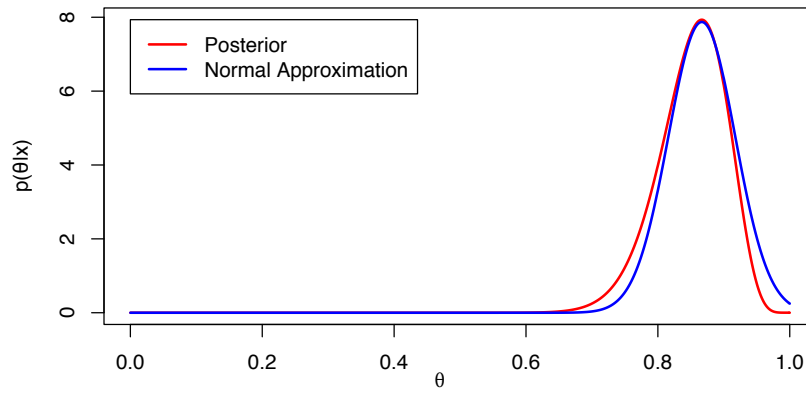


Figure 1: Posterior pdf and its Gaussian approximation for Beta(3, 3) prior on binomial success probability.

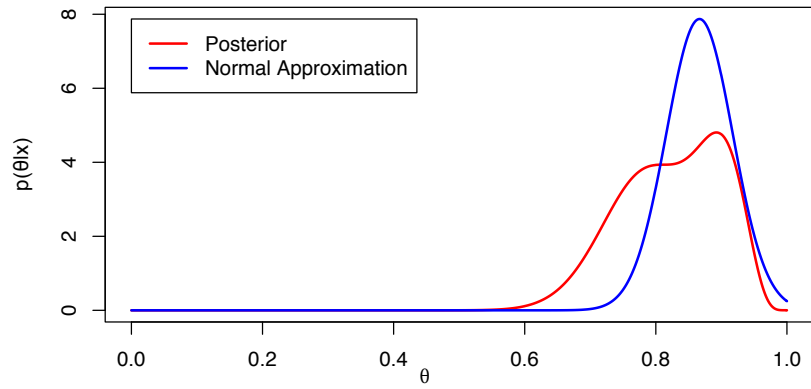


Figure 2: Posterior pdf and its Gaussian approximation for $\frac{1}{2}$ Beta(8, 2) + $\frac{1}{2}$ Beta(2, 8) prior on binomial success probability.

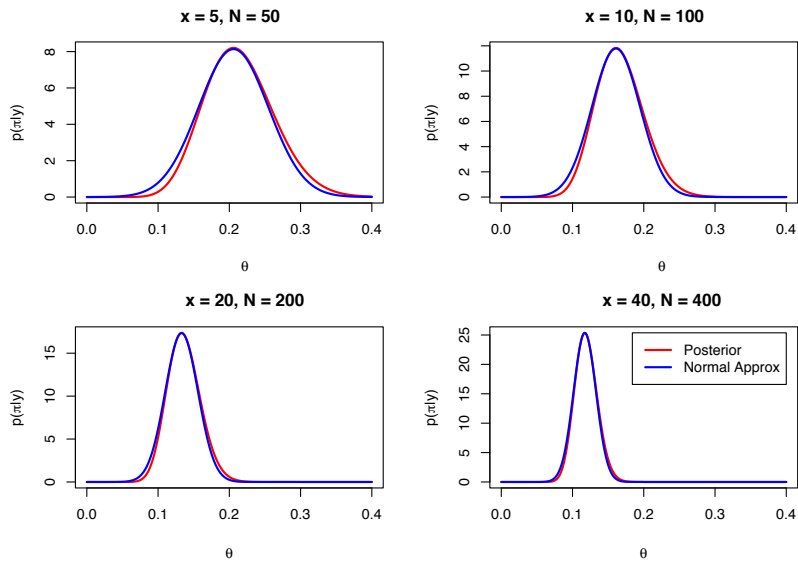


Figure 3: Posterior pdf and its Gaussian approximation for Beta(10, 10) prior on binomial success probability.

Asymptotic Normality and Consistency for Bayesian Models

CONSIDER conditionally i.i.d. observations $(X[n])_{n=0}^{N-1}$ given θ , following

$$\{X[n] \mid \Theta = \theta\} \sim f_{X|\Theta}(x \mid \theta_0) \quad (2)$$

where θ_0 is the true value of the parameter. As before, we also have

- $f_{X|\Theta}(x \mid \theta)$, data model, likelihood;
- $f_{\Theta}(\theta)$ or $p_{\Theta}(\theta)$, the prior distribution on θ .

Define $\mathbf{x} = (x[n])_{n=0}^{N-1}$. For simplicity, we consider the case of a scalar parameter θ , but the results can be generalized.

In the following discussion, we assume that the data model is correct and that θ_0 is the unique minimizer of

$$D(f_{X|\Theta}(x[n] \mid \theta_0) \parallel f_{X|\Theta}(x[n] \mid \theta)) = E_{X|\Theta} \left(\ln \frac{f_{X|\Theta}(X[n] \mid \theta_0)}{f_{X|\Theta}(X[n] \mid \theta)} \mid \theta_0 \right). \quad \text{identifiability condition}$$

Theorem 1 (Convergence in discrete parameter space). *If the parameter space Θ is finite and*

$$p_{\Theta}(\theta_0) = \Pr_{\Theta} \{\Theta = \theta_0\} > 0$$

then

$$p_{\Theta|\mathbf{X}}(\theta_0 \mid \mathbf{x}) = \Pr_{\Theta|\mathbf{X}} \{\Theta = \theta_0 \mid \mathbf{X} = \mathbf{x}\} \rightarrow 1 \quad \text{as } N \nearrow +\infty.$$

Proof. Consider the log posterior odds:

$$\ln \left[\frac{p_{\Theta|\mathbf{X}}(\theta \mid \mathbf{x})}{p_{\Theta|\mathbf{X}}(\theta_0 \mid \mathbf{x})} \right] = \ln \left[\frac{p_{\Theta}(\theta)}{p_{\Theta}(\theta_0)} \right] + \sum_{n=0}^{N-1} \ln \left[\frac{f_{X|\Theta}(x[n] \mid \theta)}{f_{X|\Theta}(x[n] \mid \theta_0)} \right]. \quad (3)$$

The second term in this expression is the sum of N conditionally i.i.d. random variables given θ_0 . Recall that $(X[n])_{n=0}^{N-1}$ are coming from $f_{X|\Theta}(x[n] \mid \theta_0)$; then

$$E_{X|\Theta} \left[\ln \left(\frac{f_{X|\Theta}(X[n] \mid \theta)}{f_{X|\Theta}(X[n] \mid \theta_0)} \right) \mid \theta_0 \right] = -D(f_{X|\Theta}(x[n] \mid \theta_0) \parallel f_{X|\Theta}(x[n] \mid \theta)) \leq 0$$

where, by our assumption that θ_0 is the unique minimizer of $D(f_{X|\Theta}(X \mid \theta_0) \parallel f_{X|\Theta}(X \mid \theta))$, the equality holds only for $\theta = \theta_0$.

If $\theta \neq \theta_0$, the second term in (3) is the sum of N i.i.d. random variables with negative mean, which diverges to $-\infty$ as $N \nearrow \infty$. As long as

$$\Pr_{\Theta}(\Theta = \theta_0) = p_{\Theta}(\theta_0) > 0$$

Terminology. We refer to the posterior ratio

$$\frac{p_{\Theta|\mathbf{X}}(\theta \mid \mathbf{x})}{p_{\Theta|\mathbf{X}}(\theta_0 \mid \mathbf{x})}$$

as **posterior-odds ratio**: θ versus θ_0 , in this case.

then the first summand in (3) is finite

the log posterior odds in (3) $\searrow -\infty$ as $N \nearrow +\infty$. Thus, if $\theta \neq \theta_0$, the posterior odds go to zero:

$$\frac{p_{\Theta|X}(\theta | \mathbf{x})}{p_{\Theta|X}(\theta_0 | \mathbf{x})} \rightarrow 0$$

which implies $p_{\Theta|X}(\theta | \mathbf{x}) \searrow 0$. As all the probabilities summed over all values of θ must add to one, we have

$$p_{\Theta|X}(\theta_0 | \mathbf{x}) \rightarrow 1.$$

□

Theorem 2 (Convergence in continuous parameter space). *If θ is defined on a compact set (i.e., closed and bounded) and A is an open subset of the parameter space containing θ_0 with prior probability $f_{\Theta}(\theta)$ satisfying $\int_{\theta \in A} f_{\Theta}(\theta) d\theta > 0$, then*

$$\Pr_{\Theta|X} \{ \Theta \in A | X = \mathbf{x} \} \rightarrow 1 \quad \text{as } N \nearrow \infty.$$

Proof. Similar to the proof for the discrete case. □

* TECHNICAL details:

- In many popular continuous-parameter scenarios, the parameter space is not a compact set, e.g. the parameter space for the mean of a Gaussian random variable is $(-\infty, +\infty)$. Luckily, for most problems of interest, the compact-set assumption of Theorem 2 can be relaxed.
- Similarly, Theorem 1 can often be extended to allow for an infinite discrete parameter space.

Theorem 3 (Asymptotic Normality of $f_{\Theta|X}(\theta | \mathbf{x})$). *Under some regularity conditions (particularly that θ_0 is not on the boundary of the parameter space) and under the conditional i.i.d. measurement model (2),*

$$\sqrt{N} \left[\hat{\theta}_{\text{MAP}}(\mathbf{X}) - \theta_0 \right] \xrightarrow{d} \mathcal{N}(0, \mathcal{I}_1^{-1}(\theta_0)) \quad \text{as } N \nearrow +\infty$$

where $\mathcal{I}_1(\theta_0)$ is the Fisher information for $\theta = \theta_0$ and a single measurement (say $X[0]$):

$$\begin{aligned} \mathcal{I}_1(\theta_0) &= \mathbb{E}_{X[0]|\Theta} \left[\left(\frac{d \ln f_{X|\Theta}(X[0]|\theta)}{d\theta} \right)^2 \middle| \theta_0 \right] \\ &= -\mathbb{E}_{X[0]|\Theta} \left[\frac{d^2 \ln f_{X|\Theta}(X[0]|\theta)}{d\theta^2} \middle| \theta_0 \right]. \end{aligned}$$

Proof. See [Gelman et al. 2014, App. B]. □

Here are some useful observations to help justify Theorem 3. Consider the scalar version of the Taylor-series expansion (1):

$$\ln f_{\Theta|\mathbf{X}}(\theta | \mathbf{x}) \approx \ln f_{\Theta|\mathbf{X}}(\hat{\theta}_{\text{MAP}} | \mathbf{x}) + 0.5(\theta - \hat{\theta}_{\text{MAP}})^2 \frac{d^2 \ln f_{\Theta|\mathbf{X}}(\hat{\theta}_{\text{MAP}} | \mathbf{x})}{d\theta^2}.$$

Now, study the behavior of the negative Hessian of the log posterior pdf at θ_0 :

$$\begin{aligned} -\frac{d^2 \ln f_{\Theta|\mathbf{X}}(\theta_0 | \mathbf{x})}{d\theta^2} &= -\frac{d^2 \ln f_{\Theta}(\theta_0)}{d\theta^2} - \frac{d^2 \ln f_{\mathbf{X}|\Theta}(\mathbf{x} | \theta_0)}{d\theta^2} \\ &= -\frac{d^2 \ln f_{\Theta}(\theta_0)}{d\theta^2} - \sum_{n=0}^{N-1} \frac{d^2 \ln f_{\mathbf{X}|\Theta}(x[n] | \theta_0)}{d\theta^2} \end{aligned}$$

and, therefore,

$$\mathbb{E}_{\mathbf{X}|\Theta} \left[-\frac{d^2 \ln f_{\Theta|\mathbf{X}}(\theta_0 | \mathbf{x})}{d\theta^2} \middle| \theta_0 \right] = -\frac{d^2 \ln f_{\Theta}(\theta_0)}{d\theta^2} + N\mathcal{I}_1(\theta_0)$$

implying that, as N grows,

$$-\frac{d^2 \ln f_{\Theta|\mathbf{X}}(\theta_0 | \mathbf{x})}{d\theta^2} \approx N\mathcal{I}_1(\theta_0).$$

To summarize: For a large number of i.i.d. measurements (i.e., asymptotically), the MAP and maximum-likelihood (ML) estimates give equivalent answers.

References

Gelman, Andrew, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin (2014). *Bayesian Data Analysis*. 3rd ed. Boca Raton, FL: Taylor & Francis (cit. on pp. 1, 5).