

Bayesian inference (continued)

Rasmus Waagepetersen
Department of Mathematics
Aalborg University
Denmark

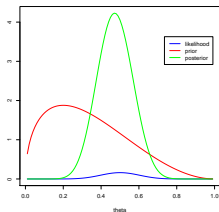
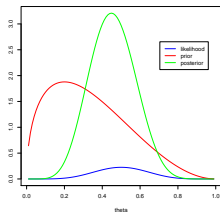
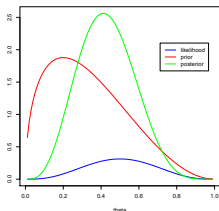
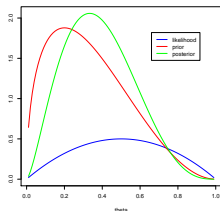
April 13, 2021

Today: selected topics in Bayesian statistics.

Asymptotics, REML, improper Gaussian

Influence of number of observations/convergence of posterior - binomial-beta

Beta-prior ($\alpha = 1.5$ $\beta = 3$). Observations $x = n/2$,
 $n = 2, 6, 12, 24$. Posterior mode 0.33, 0.41, 0.44, 0.47.



Note: posterior appears to converge to a normal density !

Bayesian asymptotics

Consider posterior of θ given observations y_1, \dots, y_n . Let $\hat{\theta}_n$ and $i_n(\theta)$ denote the MLE and Fisher information based on y_1, \dots, y_n .

Under appropriate regularity conditions, as $n \rightarrow \infty$,

$$\sup_A |P(i_n(\hat{\theta}_n)^{-1/2}(\theta - \hat{\theta}_n) \in A | y_1, \dots, y_n) - P(Z \in A)| \rightarrow 0$$

where $Z \sim N(0, I)$.

That is, posterior distribution of $i_n(\hat{\theta}_n)^{-1/2}(\theta - \hat{\theta}_n)$ converges in *total variation* distance (and hence in distribution) to the standard normal distribution. Note: given y_1, \dots, y_n , $\hat{\theta}_n$ is fixed !

Standard frequentist theory gives $i_n(\hat{\theta})^{-1/2}(\hat{\theta}_n - \theta)$ converges in distribution to a standard normal distribution but in this case θ represents fixed 'true' value while randomness of $\hat{\theta}_n$ due to sampling variation.

REML as marginal likelihood

Bayesian derivation of REML.

Consider linear mixed model $Y \sim N(X\beta, V(\psi))$.

Assume improper prior $p(\beta|\psi) \propto 1$.

Then REML is obtained by integrating out β in 'joint density' of (Y, β) :

$$\text{REML} = f(y; \psi) = \int f(y|\beta, \psi)p(\beta|\psi)d\beta$$

To show this we first compute $f(y; \psi)$ and compare it with REML.

Let $V(\psi) = LL^T$ and $\tilde{Y} = L^{-1}Y$. Then $\tilde{Y}|\beta \sim N(\tilde{X}\beta, I)$ where $\tilde{X} = L^{-1}X$. Moreover (applying Pythagoras),

$$f(\tilde{y}|\psi) = \int f(\tilde{y}|\beta, \psi) d\beta = (2\pi)^{(p-n)/2} |\tilde{X}^T \tilde{X}|^{-1/2} \exp\left(-\frac{1}{2} \|\tilde{Y} - \tilde{X}\hat{\beta}\|^2\right)$$

where $\hat{\beta}$ is the MLE. Thus (using the transformation theorem)

$$f(y|\psi) = \frac{(2\pi)^{(p-n)/2}}{|V(\psi)|^{1/2} |X^T V^{-1}(\psi) X|^{1/2}} \exp\left[-\frac{1}{2} \tilde{Y}^T (I - \tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T) \tilde{Y}\right]$$

Moreover,

$$\tilde{Y}^T (I - \tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T) \tilde{Y} = \|\tilde{Y} - \tilde{X}\hat{\beta}\|^2 = (Y - X\hat{\beta})^T V^{-1} (Y - X\hat{\beta})$$

(which may explain why REML can be short for “residual MLE”)

REML is likelihood of $A^T Y \sim N(0, A^T V(\psi)A)$. Let $\tilde{A} = L^T A$.
Then $\tilde{A}^T \tilde{X} = 0$ and

$$\tilde{A}(\tilde{A}^T \tilde{A})^{-1} \tilde{A}^T = I - \tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T.$$

Thus

$$\begin{aligned} \text{REML} &= |A^T V(\psi)A|^{-1/2} \exp\left[-\frac{1}{2} Y^T A(A^T V(\psi)A)^{-1} A^T Y\right] = \\ &|A^T V(\psi)A|^{-1/2} \exp\left[-\frac{1}{2} \tilde{Y}^T (I - \tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T) \tilde{Y}\right] \end{aligned}$$

We hence just need to show that

$$|A^T V(\psi)A| = \text{const} |V(\psi)| |X^T V^{-1}(\psi)X|$$

This follows from

$$\begin{aligned} |A^T A| |X^T X| |V| &= \left| \begin{bmatrix} A^T A & 0 \\ 0 & X^T X \end{bmatrix} V \right| = \left| \begin{bmatrix} A^T \\ X^T \end{bmatrix} [A \ X] V \right| = \\ & \left| \begin{bmatrix} A^T \\ X^T \end{bmatrix} V [A \ X] \right| = \left| \begin{bmatrix} A^T V A & A^T V X \\ X^T V A & X^T V X \end{bmatrix} \right| = \\ & |A^T V A| |X^T V X - X^T V A (A^T V A)^{-1} A^T V X| = \\ & |A^T V A| |X^T X (X^T V^{-1} X)^{-1} X^T X| = |A^T V A| |X^T X|^2 |X^T V^{-1} X|^{-1} \end{aligned}$$

(recall for partitioned matrix B , $|B| = |B_{11}| |B_{22} - B_{21} B_{11}^{-1} B_{12}|$ and $A(A^T V A)^{-1} A^T = V^{-1} - V^{-1} X (X^T V^{-1} X)^{-1} X^T V^{-1}$)

Thus

$$|A^T V(\psi) A| = \frac{|A^T A|}{|X^T X|} |V(\psi)| |X^T V^{-1}(\psi) X|$$

Pairwise difference prior

Suppose we observe $Y_i \sim N(\theta_i, 1)$, $i = 1, \dots, n$ and we want to infer $\theta_1, \dots, \theta_n$. Prior information: θ_i and θ_{i+1} “similar”.

Consider stationary AR(1) prior ($\tau_1^2 = \tau^2/(1-a)$):

$$f(\theta; a) \propto \exp\left(-\frac{1}{\tau_1^2}\theta_1^2\right) \prod_{i=2}^n \exp\left[-\frac{1}{2\tau^2}(\theta_i - a\theta_{i-1})^2\right]$$

Consider $a \rightarrow 1$. Then “limit” of right hand side is

$$f(\theta) \propto \exp\left[-\frac{1}{2\tau^2} \sum_{i=2}^n (\theta_i - \theta_{i-1})^2\right] = \exp\left[-\frac{1}{2\tau^2} \theta^T Q \theta\right]$$

for which

$$Q = \begin{bmatrix} 1 & -1 & 0 & \dots \\ -1 & 2 & -1 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \dots & -1 & 2 & -1 \\ \dots & 0 & -1 & 1 \end{bmatrix}$$

Note $Q\mathbf{1}_n = 0$ so Q does not have full rank.

On the other hand $Qx = 0$ implies $x = a1_n$ for some $a \in \mathbb{R}$. This follows since $Q = D^T D$ where D is the $(n-1) \times n$ matrix

$$D = \begin{bmatrix} -1 & 1 & 0 & \dots \\ 0 & -1 & 1 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \dots & 0 & -1 & 1 \end{bmatrix}$$

Hence the null space N_Q of Q is the span of 1_n . Note 1_n is the last eigenvector of Q with eigenvalue 0.

$f(\theta)$ not a proper density on \mathbb{R}^n since Q is not positive definite.

Limiting posterior may nevertheless still be proper (exercise).

$f(\theta)$ is invariant to addition of a constant to all elements of θ (filters constants). Hence does not imply prior assumptions about 'level' of data $Y_i | \theta_i \sim N(\theta_i, 1)$. Only need to choose prior parameter τ^2 (smoothness parameter).

Prediction

Suppose we want to predict Y_2 given Y_1 where both depend on θ .

Frequentist approach: use $f(y_2|y_1, \hat{\theta})$ where $\hat{\theta}$ estimate based on y_1 . This in general ignores extra uncertainty due to replacing θ by an estimate.

Bayesian approach offers a systematic way to take into account uncertainty of parameters in prediction by integrating out unknown parameters:

$$f(y_2|y_1) = \int f(y_2|y_1, \theta)p(\theta|y_1)d\theta = \mathbb{E}_{\theta|y_1} f(y_2|y_1, \theta)$$

Predictive density is 'weighted' average of predictive densities $f(y_2|y_1, \theta)$ where 'weights' given by posterior density $p(\theta|y_1)$ reflects uncertainty of θ .

Nice solution in principle but in practice the computation may not be straightforward.

Difficult posteriors

Except for the simple examples with conjugate priors the posterior is often intractable - closed form expressions for posterior quantities like expectations, variances, quantiles etc. often not available.

Non-normal example: logistic regression with normal prior

$\beta \sim N(0, \tau^2)$ (normal prior)

$Y_j|\beta \sim \text{binomial}(n_j, p_j)$ conditionally independent given β $j = 1, \dots, n_j$

$\log(p_j/(1 - p_j)) = \eta_j = x_j^T \beta$

$p_j = \exp(\eta_j)/(1 + \exp(\eta_j))$

Likelihood function:

$$f(y|\beta) = \prod_j p_j^{y_j} (1 - p_j)^{1-y_j} = \prod_j \frac{\exp(x_j^T \beta)^{y_j}}{(1 + \exp(x_j^T \beta))^{n_j}}$$

Marginal density $f(y)$:

$$\int_{\mathbb{R}} f(y|\beta) f(\beta; \tau^2) d\beta = \int_{\mathbb{R}} \prod_j \frac{\exp(x_j^T \beta)^{y_j}}{(1 + \exp(x_j^T \beta))^{n_j}} \frac{\exp(-\beta^2/(2\tau^2))}{\sqrt{2\pi\tau^2}} d\beta$$

Integral can not be evaluated in closed form.

Laplace/Gaussian approximation

Let $g(\beta) = \log(f(y|\beta)f(\beta))$ and choose $\hat{\beta}$ so $g'(\hat{\beta}) = 0$
($\hat{\beta} = \arg \max g(\beta)$).

Note: $\hat{\beta}$ is MAP (maximum a posteriori) estimate. Not MLE.

Taylor expansion around $\hat{\beta}$:

$$g(\beta) \approx \tilde{g}(\beta) =$$

$$g(\hat{\beta}) + (\beta - \hat{\beta})g'(\hat{\beta}) + \frac{1}{2}(\beta - \hat{\beta})^2 g''(\hat{\beta}) = g(\hat{\beta}) - \frac{1}{2}(\beta - \hat{\beta})^2 (-g''(\hat{\beta}))$$

i.e. $\exp(\tilde{g}(\beta))$ proportional to normal density $N(\mu_{LP}, \sigma_{LP}^2)$,
 $\mu_{LP} = \hat{\beta}$ $\sigma_{LP}^2 = -1/g''(\hat{\beta})$.

Since

$$p(\beta|y) \approx \exp(g(\beta)) \approx \exp(\tilde{g}(\beta))$$

it follows

$$\beta|Y = y \approx N(\hat{\beta}, -1/g''(\hat{\beta}))$$

Regarding marginal density of y :

$$\begin{aligned} f(y) &= \int_{\mathbb{R}} \exp(g(\beta)) d\beta \approx \int_{\mathbb{R}} \exp(\tilde{g}(\beta)) d\beta \\ &= \exp(g(\hat{\beta})) \int_{\mathbb{R}} \exp\left(-\frac{1}{2\sigma_{LP}^2}(\beta - \mu_{LP})^2\right) d\beta = \exp(g(\hat{\beta})) \sqrt{2\pi\sigma_{LP}^2} \end{aligned}$$

Note: these kinds of arguments basis of asymptotic results for posterior distributions.

Other approaches

Numerical integration (Gaussian quadrature), Monte Carlo, importance sampling, Markov chain Monte Carlo,....

Enough material for a whole course.

Why/when is Bayesian inference useful

- ▶ obvious if prior information is available
- ▶ for highly complex models maximum likelihood inference is difficult (multimodality, evaluation of likelihood).
Computation of posterior expectations and probabilities numerically more simple.
- ▶ can compute posterior distributions of complicated parameters whose distribution may be hard to obtain in the MLE setting.
- ▶ natural approach to take into account parameter uncertainty in prediction.

Exercises

1. (hidden AR(1) model) assume that Y_i given θ are independent $N(\theta_i, 1)$ and that $\theta = (\theta_1, \dots, \theta_n)$ follows a stationary AR(1) process prior with known autoregression parameter a and noise variance τ^2 .
 - 1.1 Compute the posterior distribution of θ (e.g. use the previous results for conditional distributions in general linear mixed models).
 - 1.2 What is the limiting posterior when $a \rightarrow 1$?
 - 1.3 is the limiting prior proper ? is the limiting posterior proper ?

2. Consider data X_1, \dots, X_n from a zero-mean AR(1) process. Consider the *conditional likelihood* of X_2, \dots, X_n given (X_1, a, τ^2) .
- 2.1 Show that the posterior distribution of (a, τ^2) obtained by combining the conditional likelihood with the (improper) prior $p(a, \tau^2) \propto 1/\tau^2$ is equivalent to the posterior for a linear normal model with observation vector $(X_2, \dots, X_n)^T$ and design matrix given by the column $(X_1, \dots, X_{n-1})^T$.
- 2.2 use the previous results to compute the predictive mean and variance of X_{n+1} given X_1, \dots, X_n (again using the conditional likelihood instead of the usual likelihood of (X_1, \dots, X_n)).

NB: rather than using the conditional likelihood we could instead assume $X_1 \sim N(\mu_1, 1)$ and use the usual likelihood of (X_1, \dots, X_n) given (μ_1, a, τ^2) combined with the prior $p(\mu_1, a, \tau^2) \propto 1/\tau^2$. This would give the same posterior inference for (a, τ^2) as by using the conditional likelihood.