

Bayesian model comparison (continued)

STAT 206
4 Mar 21

Lecture

Let M be either (a) a frequentist model (with no prior layer in the hierarchy) or (b) the sampling model [SM] part of a Bayesian model (ignoring the prior).
~~(frequentist)~~

The deviance of M , as defined in the 3 Mar 21 discussion section notes, is

$$D_F(M | \mathbb{D} \mathcal{B}) = -2 \cdot \log(\hat{\theta}_{MLE} | \mathbb{D} [SM] \mathcal{B})$$

Annotations:
 - $\hat{\theta}$ is the model's parameter vector
 - \mathbb{D} is the dataset

the log likelihood function is defined by taking the log of {the likelihood function obtained directly from the SM, including the SM's normalizing constant}.

Notice

that deviance_F($M | \mathbb{D} \mathcal{B}$) is a single value

The Bayesian way of looking at the deviance ⁽²⁾
is different: $D(\theta | M, D, \mathcal{B}) = \underbrace{\text{the usual}}_{\text{likelihood function}}$

In other words, ^{deterministic}
 D_B is a function of the parameter vector θ

(D_F is just the maximum value attained by
 D_B over the parameters)

This means that

$D_B(\theta | M, D, \mathcal{B})$ has a (marginal) ^{posterior} distribution

of its own (!), which we can monitor

as a derived quantity in our (MCMC) MCMC (!)

simple
example

$(Y_i | \theta, \mathcal{B}) \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\theta), (i=1, \dots, n)$
 $\theta \in (0, 1) \quad Y = (Y_1, \dots, Y_n)$

$P(Y | \theta, \mathcal{B}) = \theta^s (1-\theta)^{n-s}$ (note: in this case
the normalizing constant is 1)

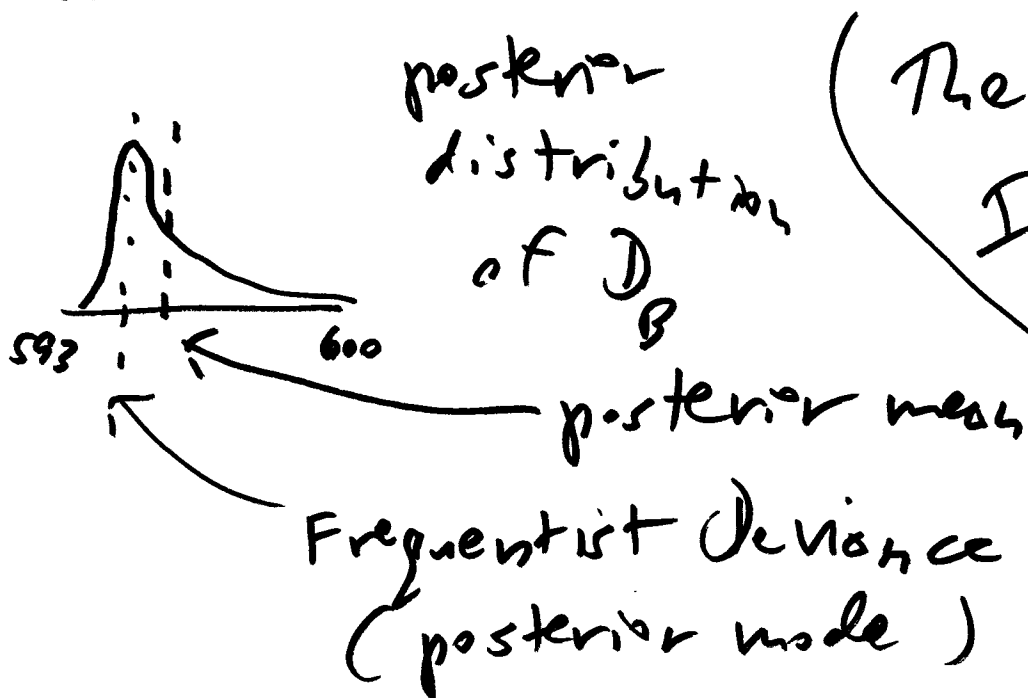
$P(\theta | Y, \mathcal{B}) = \theta^s (1-\theta)^{n-s}$

$$l(\theta | \mathbf{z} | \mathcal{B}) = s \log \theta + (n-s) \log (1-\theta) \quad (3)$$

$$D_B(\theta | M | \mathcal{B}) = -2 l(\theta | \mathbf{z} | \mathcal{B})$$

$$= -2 [s \log \theta + (n-s) \log (1-\theta)]$$

(R code)



The Deviance
Information

Criterion
(DIC)

was developed
by David
Spiegelhalter

including Nicky East & Angelika
van der Linde by David Spiegelhalter
& colleagues, in the early 2000s & published
in 2002

As described in the document

camera notes from the morning discussion
section on 3 Nov 21, the deviance itself

uses the entire dataset to evaluate the model's quality, which encourages over-

fitting, so it can't serve by itself as a good model comparison tool; it needs

to be penalized by a term that discourages overly-complex models.

So a good model comparison method looks like this: $f(\text{model fit}, \text{model complexity})$

The simplest such $f(\cdot, \cdot)$ would be additive:

(Can use the deviance for this) (what should we use for this?)

$f(\text{model fit}, \text{model complexity})$
= Deviance + (model complexity), but

↓ want small

This needs measuring (model complexity) on the same scale as the Deviance;

to get apples to apples

how should this be accomplished?

it turns out that Laplace approximations can help.

First, a note on MCSE values

in MCMC If a column in the MCMC data set behaves like an

AR(1), then

$MCSE(\bar{\mu}^*) = \frac{S_{\mu^*}}{\sqrt{M}}$

as $\rho \uparrow +1$ MCSE $\rightarrow +\infty$

$\frac{1 + \rho^2}{1 - \rho^2}$



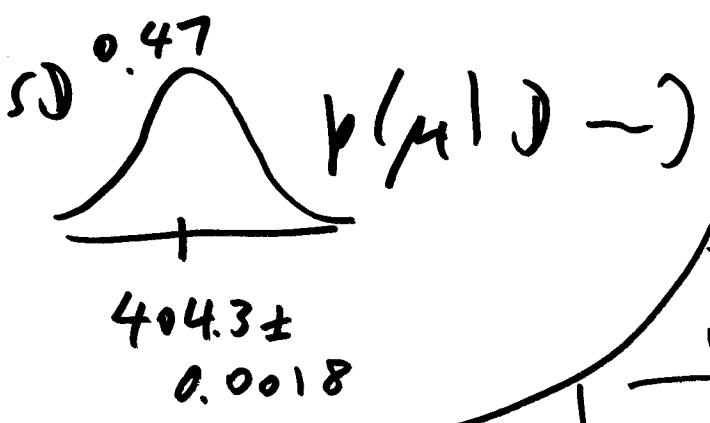
MC est. of post. mean for μ

the MCMC algorithm works by finding a Markov chain whose equilibrium (stationary) distribution is $p(\theta | D)$

Markov chain

$$\theta_{t+2}^* = f(\theta_{t+1}^*)$$

$$B+1 \begin{bmatrix} \theta_{B+1}^* \\ \vdots \\ \theta_{B+2}^* \\ \vdots \end{bmatrix}$$



NB10 case study

Using the cheating approach, we identified 2 interesting sampling models.

• Gaussian (M_1)	Q_1	Is M_2 better than M_1 ?
• t (M_2)	$A_1 \rightarrow Q_2$	
	Pursuing A_2	Better for what purpose?

turns model comparison into a decision problem (Bayesian decision theory) utility

	θ	$D(\theta \sim)$
1	θ_1^*	$-2 [S \log \theta_1^* + (h-S) \log (1 - \theta_1^*)]$
2	θ_2^*	
⋮	⋮	
⋮	⋮	
M	θ_M^*	

MC

