

Bayesian model comparison with
log scores

$$D = \mathcal{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad y_i \in \mathbb{R}$$

Lecture

$$M_j: \begin{cases} (\tilde{\theta}_j | [PM] \mathcal{B}) \sim p(\tilde{\theta}_j | [PM] \mathcal{B}) \\ (\mathcal{Y}_i | [SM: M_j] \tilde{\theta}_j, \mathcal{B}) \stackrel{i.i.d.}{\sim} p(y_i | [SM: M_j] \tilde{\theta}_j, \mathcal{B}) \end{cases} \quad (i=1, \dots, n) \quad (1)$$

As noted

$$\theta_j \in \mathbb{R}^{k_j} \quad \text{PDF}$$

last time, log scores derive from the

Prediction
Principle

Good models make good predictions,
Bad ~~models~~ ~~make~~ bad ~~predictions~~.
That's how we know a model

is good or bad \oplus of data not used in

the model-fitting process $\left\{ \begin{array}{l} \text{Bayesian prediction} \end{array} \right.$

of a new data value in the setting above

involves calculating the posterior predictive

distribution $p(y^* | \mathcal{Y}, M_j, \mathcal{B})$ for a new $\left. \begin{array}{l} \text{data} \\ \text{value} \end{array} \right\} y^*$

$$p(y^* | z^M; \mathcal{B}) = \int_{\Theta_j} p(y^*, \theta_j | z^M; \mathcal{B}) d\theta_j$$

(PDF) \nearrow (H) ↓

$$= \int_{\Theta_j} p(y^* | \theta_j; z^M; \mathcal{B}) \cdot \underline{p(\theta_j | z^M; \mathcal{B})} d\theta_j$$

$$= \int_{\Theta_j} p(y^* | \theta_j; z^M; \mathcal{B}) p(\theta_j | z^M; \mathcal{B}) d\theta_j$$

(Sampling PDF for a new data value)

(posterior distrib. for θ_j given the data z^M & \mathcal{B}) (PDF)

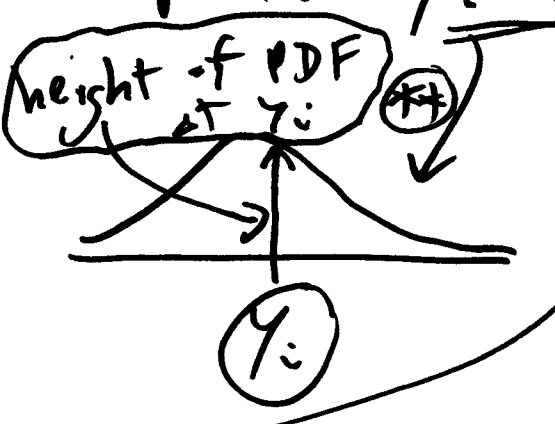
This looks formidable but is actually easy to simulate from in MCMC (see below)

OK, suppose that we have

$p(y^* | z^M; \mathcal{B})$; { how can we use it in model comparison? }

① Let's imagine jack-knifing the data set ^③ to get (data not used in the fitting process): define $\underline{y}_{[-i]} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{i-1} \\ y_{i+1} \\ \vdots \\ y_n \end{pmatrix}$ i.e. \underline{y} with y_i set aside

② Now locate y_i in $p(y_i | \underline{y}_{[-i]}, \mu, \beta)$:



ⓐ: How to make a good comparison between y_i (a scalar) and

$p(y_i | \underline{y}_{[-i]}, \mu, \beta)$ (a PDF)?

ⓐ: ~~this is~~ This is

called finding a good scoring rule; under reasonable assumptions the optimal choice is $\log [p(y_i | \underline{y}_{[-i]}, \mu, \beta)] = \log (** \text{ above})$

③ Now we compute $LS_{CV}(M_j | Z, B) =$

In my research this is called

$$\frac{1}{n} \sum_{i=1}^n \log [p(y_i | y_{[-i]}, M_j, B)]$$

the cross-validated (CV) log score

④ However, (my) research has shown the intuitively-reasonable fact that even with n only moderate in size

$$p(y_i^* | y_{[-i]}, M_j, B) \approx p(y_i^* | y, M_j, B)$$

for all y^*

this motivates

$$LS_{FS}(M_j | Z, B) = \frac{1}{n} \sum_{i=1}^n \log [p(y_i | y, M_j, B)]$$

full-sample (FS) predictive distribution for a new data point y^*

using MCMC to calculate $LS_{FS}(n_j | \gamma^*)$ (5)

From above

$$p(\gamma^* | \gamma n_j; \mathcal{B}) = \int_{\theta_j} p(\gamma^* | \theta_j; n_j; \mathcal{B}) \cdot \underbrace{p(\theta_j | \gamma n_j; \mathcal{B})}_{\text{PDF}} d\theta_j$$

This is a mixture representation of $p(\gamma^* | \gamma n_j; \mathcal{B})$, so (it) can be approximated to arbitrary MC accuracy as follows:

- (A) obtain the MCMC data set for θ_j as usual, with n (large) monitoring iterations
- (B) monitor the deterministic quantity $p(\gamma^* | \theta_j^*; n_j; \mathcal{B})$
- (C) Take the mean of the quantities in (B):

$$\underline{p(y^* | \gamma, n; \mathcal{B})} = \frac{1}{M} \sum_{m=1}^M \underline{p(y^* | \begin{pmatrix} \theta_j^* \\ \sim \end{pmatrix}_m, n; \mathcal{B})}$$

and now

$$LS_{FS}(n; \gamma | \mathcal{B}) =$$

it's as easy as
 this: ~~large LS values~~ ^{a good model}

$$\frac{1}{n} \sum_{i=1}^n \log \left[\frac{1}{M} \sum_{m=1}^M p(\underbrace{y_i}_{\sim} | \theta_j^* | n; \mathcal{B}) \right]$$

(Rcode) CRAN seems not to ~~offer~~ offer

much support for log scores, but

(a) LS_{cv} is also called the conditional

predictive ordinate (CPO) method &

(b) there are some packages that
 implement CPO

The LS_{CV} idea goes back to Geisser (5)
& Eddy (1979) & Gelfand (1996)

Q: where are (GOF) & (model complexity)

in LS_{FS} ?

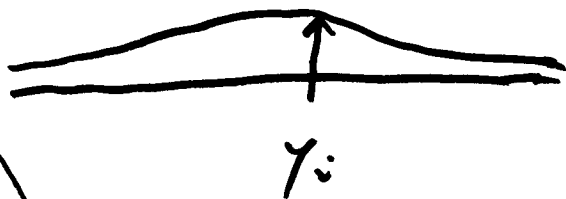
A:

$p(y_i | \gamma, M, B)$



γ_i

$p(y_i | \gamma, M_2, B)$



γ_i

needlessly spread out

(worse LS)

Note:

LS_{FS} works

with both

fixed effects

& random

effects

models

M_2 has extra parameters in it that are not helpful in prediction of future data

$$p(y_i | \underline{\mu}, \sigma)$$



γ_i

$$p(y_i | \underline{\mu}, \sigma)$$



γ_i

has extra parameters that (greatly) improve the model fit

μ_1 fit to

aspirin data:

random effects,

$$(\theta_i | \mu, \sigma^2) \stackrel{\text{IID}}{\sim} \mathcal{N}(\mu, \sigma^2)$$

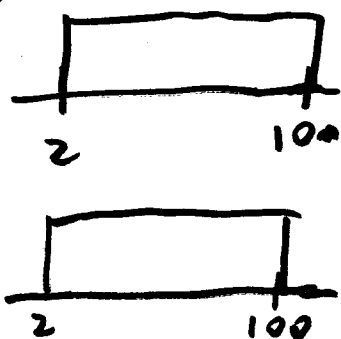
μ_2 fit —

$k=6$

ie., γ cannot be learned with



this data set



$$(\theta_i | \mu, \sigma^2) \stackrel{\text{IID}}{\sim} t_r(\mu, \sigma^2)$$

$(i=1, \dots, k)$

$p(\gamma)$

↓ same

$p(\gamma | D \dots)$

as $r \rightarrow \infty$
 $t_r \rightarrow \mathcal{N}$

Comparison of Bayesian model comparison methods

BIC, other Bayes factors

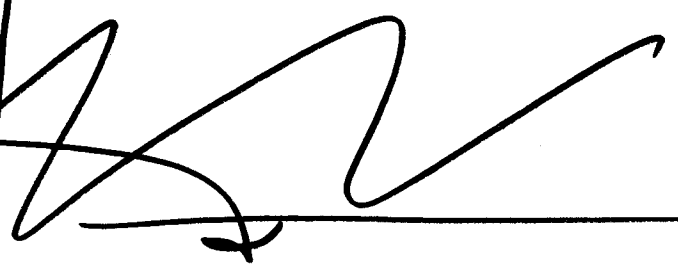
DIC, AIC, LS

designed with what goal in mind?



achieve frequentist consistency at both M_1 & M_2
 (more complex)
 complex

BIC, other Bayes factors tend to favor less complex models than {DIC, AIC, LS}



favor models that make good predictions
 0.05 out of 100 sample

achieve frequentist consistency at M_2 but not at M_1

models that predict 0.05 well tend to be (somewhat) more complex than models that achieve

frequentist consistency of Bayesian ⁽¹⁰⁾
model comparison methods \leftrightarrow vantage
to be a well-calibrated Bayesian

~~***~~ assumptions / ^(more complex)
① either $M_1 = M_{DG}$ or $M_2 = M_{DG}$
② $n \rightarrow \infty$ holding M_1 & M_2 fixed

$$P_F(\text{BIC chooses } M_1 \mid M_1 = M_{DG}) \xrightarrow{n \rightarrow \infty} 1$$

if so, BIC / ^{Bayes factors} ~~is~~ are consistent at the
less complex model : true

$$P_F(\text{BIC chooses } M_2 \mid M_2 = M_{DG}) \xrightarrow{n \rightarrow \infty} 1$$

if so BIC / Bayes factors are consistent
at the more complex model : true

However (LS, DIC, AIC) (1)

$$P_F(\text{LS chooses } M_2 \mid M_2 = M_{DG}) \xrightarrow{n \rightarrow \infty} 1$$

True

But

$$P_F(\text{LS chooses } M_1 \mid M_1 = M_{DG}) \xrightarrow{n \rightarrow \infty} 0 \quad (!)$$

However,

Q: M_1 $n = 79$; how do these asymptotic results apply, if at all?

A: not always

Q: " $n \rightarrow \infty$ holding k_1 & k_2 fixed "

is not realistic for actual data science.

as $n \uparrow$, to decently model real-world

complexity need $k_1 \uparrow$ & $k_2 \uparrow$ better asymptotics: $k_0, n \rightarrow \infty$ with $0 < \frac{k}{n} \rightarrow c < 1$