

Laplace

In the context of Bayes

STAT 206

5 Mar 21

Approximations factors as an approach

DD AM (5 Mar) discussion section

to Bayesian model comparison

9 Mar: lecture

Remember Bayes factors from week 1?

Suppose that we have uncovered an ensemble  $\mathcal{M} = \{M_1, M_2, \dots, M_m\}$  of worthwhile-to-consider models (e.g., using the cheating

approach)

$\underline{P} = (\underline{Q}, \underline{C}) + (\underline{\theta}, \underline{D}, \underline{B})$  ← unknown quantities of principal interest

→  $M_j$  :  $(j=1, \dots, m)$

$(\underline{\theta}_j | [PM_j] \underline{B}) \sim p(\underline{\theta}_j | [PM_j] \underline{B})$   
 $(\underline{Y}_i | \underline{\theta}_j [SM_j] \underline{B}) \stackrel{i.i.d.}{\sim}$

$\underline{D} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$

$(\underline{Y}_i | \underline{\theta}_j [SM_j] \underline{B}) \sim p(Y_i | \underline{\theta}_j [SM_j] \underline{B})$   
 $(k_j = \# \text{ free parameters in } \underline{\theta}_j)$

ex. NB10 care study)  $\mathcal{M} = \{ \text{Gaussian, } t \}$   
 $(n=2)$

$M_1 : \begin{cases} \underline{\theta}_1 = (\mu, \sigma) \\ (\underline{\theta}_1 | [PM_1] \underline{B}) \sim p(\underline{\theta}_1 | [PM_1] \underline{B}) \\ (\underline{Y}_i | \underline{\theta}_1 [SM_1: N] \underline{B}) \stackrel{i.i.d.}{\sim} p(Y_i | \underline{\theta}_1 [SM_1: N] \underline{B}) \end{cases}$

$$M_2: \left\{ \begin{array}{l} (\theta_2 | [PM_2] B) \sim p(\theta_2 | [PM_2] B) \\ (Y_i | \theta_2 [SM_2: T] B) \sim p(Y_i | \theta_2 [SM_2: T] B) \end{array} \right\} \quad (2)$$

$(i=1, \dots, n)$   $\theta_2 = (\mu, \sigma_t, \tau)$  not the same as  $\sigma_N$  in  $M_1$

$\theta = \mu$  has the same meaning in both models

Since  $|M_2| = m < \infty$  (i.e., the number of models in  $M$  is finite) it's sufficient to have a method that compares models 2 at a time:  $M_1$  vs.  $M_2$ ; better of  $(M_1, M_2)$  vs.  $M_3$ ; best  $M$  so far vs.  $M_4$ ; ...

The Bayes factor strategy makes the following provisional (& heroic) assumptions: ① There is a "true" data-

generating model  $M_{DG}$ ; ②  $M_{DG} \in \underline{M}$ .  
(this is called the  $M$ -closed viewpoint)

Then first pretend that  $\mathcal{M} = \{M_1, M_2\}$ ; then <sup>(3)</sup>

$$\left[ \frac{P(M_2 | D, B)}{P(M_1 | D, B)} \right]^* = \left[ \frac{P(M_2 | B)}{P(M_1 | B)} \right] \left[ \frac{P(D | M_2, B)}{P(D | M_1, B)} \right]$$

posterior odds  
in favor of  $M_2$   
over  $M_1$ , given  
 $D$  &  $B$

prior odds  
in favor  
of  $M_2$  over  
 $M_1$ , given  
 $B$

Bayes factor  
in favor of  
 $M_2$  over  $M_1$ ,  
given  $D$  and  
 $B$

Q: How specify prior model probabilities?

$P(M_j | B)$ !

A: "Please ask me an easier question!"

Seriously: many people try to "avoid" this issue by assigning a uniform prior across models:

$$P(M_j | B) = \begin{cases} \frac{1}{n} & \text{for } j=1, \dots, n \text{ in } \mathcal{M} = \{M_1, \dots, M_n\} \\ 0 & \text{else} \end{cases}$$

Another idea is to penalize more complex models in the prior (details later if time).  
(Occam's Razor) ④

This task is difficult: what ~~was~~<sup>were</sup> your prior odds in favor of the  $t$  model over the Gaussian model in the NB10 case study before you saw the data?

One possible

(DR)

decision rule in Bayesian model comparison:

DR:

Choose  $M_{j^*}$  such that  $P(M_{j^*} | DB)$  is the largest posterior model probability across the models in  $\mathcal{M} = \{M_1, \dots, M_m\}$ . This

requires specifying the prior model probabilities in  $\textcircled{4}$  above. One idea: use sensitivity analysis over a variety of PMPs;

if all reasonable PMP choices lead to the same best model, breathe a sigh of relief; if not, "think harder"

Many people punt on this issue by using just the Bayes factors to decide:

$DR_2$  choose the model with the highest Bayes factor in favor of it (which is equivalent to  $DR_1$  with a uniform  $(\frac{1}{m})$  prior on  $\mathcal{M} = \{M_1, \dots, M_m\}$ )

The terms in the Bayes factor are of the form  $p(D | M_j; B)$ ; how compute or approximate?

this is called the marginal likelihood or model evidence

It's actually the prior predictive distribution for the dataset  $D$  before it's been observed, given only that  $(M_j = M_{D,G}) \& B$ .

$p(D | M_j; \mathcal{B})$  is hard to think about without <sup>(6)</sup> knowing  $(\theta_j)$ , so let's bring  $(\theta)$  into the calculations in the usual way:

$$p(D | M_j; \mathcal{B}) = \int p(D | \theta, M_j; \mathcal{B}) p(\theta_j | \mathcal{B}) d\theta_j$$

But  $p(D | \theta_j; M_j; \mathcal{B})$  is just the sampling distribution for the dataset  $D$  under model  $M_j$  with parameter vector  $\theta_j$  of  $M_j$ .

~~So this says to take a weighted average of the sampling distributions in  $M_j$ , weighted by the prior distribution for the parameter vector  $\theta_j$  of  $M_j$ .~~

So this says to take a weighted average of the sampling distributions in  $M_j$ , weighted by the prior distribution for the parameter vector  $\theta_j$  of  $M_j$ .

serious problem: if we try to use a low-information prior for  $\theta_j$ ,  $p(D | M_j; \mathcal{B})$  can be extremely & uncomfortably sensitive to the LI details.

Simple example

$$(\mu | [PM: LI] \mathcal{B}) \sim N(0, 1000) \quad (7)$$

$$(y_i | \mu [SM: N] \mathcal{B}) \stackrel{IID}{\sim} p(y_i | \mu [SM: N] \mathcal{B})$$

NB10 data

Call this model  $M^*$

$$(i=1, \dots, n) \quad (k=1) \quad \theta = (\mu)$$

$$N(\mu, \sigma^2)$$

$$p(y | \mu [SM: N] \mathcal{B}) =$$

$$\prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp \left[ -\frac{1}{2\sigma^2} (y_i - \mu)^2 \right]$$

$$= \sigma^{-n} (2\pi)^{-\frac{n}{2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right]$$

(known) set  $\sigma = 5$  (sample size is NB10 data)

consider this a function of  $\mu$  for fixed  $y$  known

$y = (y_1, \dots, y_n)$  and  $\sigma$   
NB10 data vector

LI prior for inference:  $N(0, 1000)$

evidence for  $M^*$

$$p(y | M^* \mathcal{B}) = \int_{-\infty}^{\infty} \sigma^{-n} (2\pi)^{-\frac{n}{2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right] d\mu$$

what if we tried to approximate this integral with (IID) MC!

The basic (FFD) MC algorithm for approximating <sup>⑧</sup> an integral of the form 
$$(**) = \int g(\theta) \cdot \boxed{p(\theta)} d\theta$$

PDF

is simply as follows:

① draw  $\theta_j^*$  from  $p(\theta)$  for  $j=1, \dots, M$  <sup>IID</sup>

large number of MC draws

② approximate  $(**)$  by

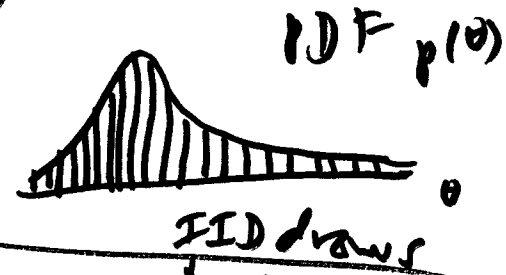
④ 
$$\int g(\theta) \cdot p(\theta) d\theta =$$

that's all there is to it.

$$\frac{1}{M} \sum_{i=1}^M g(\theta_i^*)$$

To see why.

imagine approximating  $p(\theta)$  with many equally-wide histogram bars



Sampling  $\theta_j^*$  from  $p(\theta)$  will cause each histogram bar to have a # of  $\theta_j^*$  values that's (approx.) proportional to  $p(\theta)$  (hist. bar center)

$(**)$  is a weighted average of  $g(\theta)$  values, weighted by the  $p(\theta)$  values.



(R code) Applying the MC integration idea <sup>(9)</sup>  
 to marginal likelihoods is now straightforward

$$p(D | M; \mathcal{B}) = \frac{1}{M} \sum_{i=1}^M p(D | \theta_i^* M; \mathcal{B}),$$

in which  $\theta_i^* \stackrel{\text{IID}}{\sim} p(\theta_j | \mathcal{B})$  for  $i=1, \dots, M$

(R code) But if we want  $p(\theta_j | \mathcal{B})$  to be  
 a LI prior, we must be careful in how  
 we specify such a prior (ie., LI prior <sup>(\*)</sup>)

specification in Bayesian model comparison

is different from LI prior specification <sup>(\*\*)</sup>

in Bayesian inference & prediction

& harder: why MC to approximate  
 a marginal likelihood with a LI prior  
 → next page <sup>(\*\*\*)</sup>

many LI priors  
 in inference  
 lead to the same  
 posterior

\*\*\* can (as mentioned above) be extremely sensitive to the LI prior details

I've

found that the following works well: ① do enough of a likelihood analysis to identify the interval ~~for~~ 99.99% for each parameter:

suppose that ~~only~~ this interval is (say)  $(\theta_{ij}^{left}, \theta_{ij}^{right})$

for parameter  $\theta_{ij}$  (and similarly for the other parameters) ② Specify the LI prior as follows:

$(\theta_{ij} | \mathcal{M}: \mathcal{L} \& \mathcal{B}) \sim U((\theta_{ij}^{left}, (\theta_{ij}^{right}))$   
(and so on (independently) for the other parameters)

(R code)

the Bayes factor in favor of  $M_2$  over  $M_1$  is of the form  $BF(M_2 || M_1 | \mathcal{D} \& \mathcal{B}) =$

$\frac{p(\mathcal{D} | M_2, \mathcal{B})}{p(\mathcal{D} | M_1, \mathcal{B})}$

marginal likelihood values  $p(\mathcal{D} | M_j, \mathcal{B})$  tend to be

really small (eg.,  $10^{-100}$ ), so the log function <sup>(11)</sup>  
 goes to the rescue again:

Choose  $M_2$  over  $M_1$  if  $BF(M_2 || M_1 | D, \mathcal{B}) > 1$   $\leftrightarrow$  3 equivalent decision rules

if  $p(D | M_2, \mathcal{B}) > p(D | M_1, \mathcal{B})$   $\leftrightarrow$

if  $\log p(D | M_2, \mathcal{B}) > \log p(D | M_1, \mathcal{B})$

so attention shifts to  $\log p(D | M_j, \mathcal{B})$

here's where Mr. Laplace's approximation

idea comes in:

$$\log p(D | M_j; \mathcal{B}) = \int_{\Theta_j} \underbrace{p(D | \theta_j; M_j; \mathcal{B})}_{\text{integrand}} \cdot \underbrace{p(\theta_j | M_j; \mathcal{B})}_{\text{un-normalized}} d\theta_j$$

here's Mr. Laplace's

idea: ~~the  $\theta_j$  is a random variable~~

~~the  $\theta_j$  is a random variable~~  
 $p(D | \theta_j; M_j; \mathcal{B}) =$  likelihood function for  $\theta_j$  in  $M_j$

and  $p(\theta_j | M_j, \mathcal{B})$  is the prior for  $\theta_j$  in model  $M_j$ , so the integrand is the un-normalized posterior for  $\theta_j$  under model

$M_j$  for even moderate  $n$  when  $\mathcal{D} = \mathcal{Z} = (y_1, \dots, y_n)$ ,

this posterior should be approximately multivariate normal with mean vector  $\hat{\theta}_j$  <sup>MLE</sup> and covariance matrix  $\hat{\Sigma}_{\theta_j} = \hat{I}_{\theta_j}^{-1}$ , the inverse of the observed information matrix

Let's integrate the integrand only over that subset of  $\Theta_j$  where the integrand is positive (the other parts of  $\Theta_j$  contribute zero to the integral); then we can write the integrand as  $x > 0 \rightarrow x = \exp\{\dots\}$

as  $p(\mathcal{D} | \theta_j, M_j, \mathcal{B}) p(\theta_j | M_j, \mathcal{B}) = \exp\{\log[\dots]\}$  (nice trick)

As we did in computing the deviance, let's (12)  
 choose the likelihood function  $l(\underline{\theta}_j | \mathcal{D}; \mathcal{M}_j; \mathcal{B})$   
 so that it exactly equals the sampling dist.  
 $p(\mathcal{D} | \underline{\theta}_j; \mathcal{M}_j; \mathcal{B})$  (this uniquely defines what  
 would otherwise be the arbitrary multiplicative  
 constant in the definition of the likelihood)

then integrand =  $\exp \left\{ \underbrace{l(\underline{\theta}_j | \mathcal{D}; \mathcal{M}_j; \mathcal{B})}_{s(\underline{\theta}_j)} + \log p(\underline{\theta}_j | \mathcal{M}_j; \mathcal{B}) \right\}$

Mr. Laplace now expands this function  $s(\underline{\theta}_j)$  in a  
 3-term Taylor series approximation around

the point  $\underline{\theta}_j = \begin{pmatrix} \underline{\theta}_j \\ \underline{\nu}_j \end{pmatrix}$  For simplicity let's look  
 at the special case in which  $k_j = \begin{pmatrix} \# \text{ of free} \\ \text{parameters} \end{pmatrix} = 1$   
 is a scalar in  $\mathcal{M}_j$

i.e.,  $\underline{\theta}_j = \theta_j$

The Taylor series approximation

then looks like this:

$$g(\theta_j) = g(\hat{\theta}_j) + \left. \left[ \frac{d}{d\theta_j} g(\theta_j) \right] \right|_{\theta_j = \hat{\theta}_j} \cdot (\theta_j - \hat{\theta}_j) + \frac{1}{2} \left. \left[ \frac{d^2}{d\theta_j^2} g(\theta_j) \right] \right|_{\theta_j = \hat{\theta}_j} \cdot (\theta_j - \hat{\theta}_j)^2 \quad (18)$$

substituting in gives the following:

$$\ell(\theta_j | D_n; \mathcal{B}) + \log p(\theta_j | n; \mathcal{B}) = \ell(\hat{\theta}_j | D_n; \mathcal{B}) + \log p(\hat{\theta}_j | n; \mathcal{B}) + \left\{ \left[ \frac{d}{d\theta_j} \ell(\theta | D_n; \mathcal{B}) \right]_{\theta = \hat{\theta}_j} + \left[ \frac{d}{d\theta_j} \log p(\theta | n; \mathcal{B}) \right]_{\theta = \hat{\theta}_j} \right\} \cdot (\theta - \hat{\theta}_j) + \frac{1}{2} \left\{ \left[ \frac{d^2}{d\theta_j^2} \ell(\theta | D_n; \mathcal{B}) \right]_{\theta = \hat{\theta}_j} + \left[ \frac{d^2}{d\theta_j^2} \log p(\theta | n; \mathcal{B}) \right]_{\theta = \hat{\theta}_j} \right\} \cdot (\theta_j - \hat{\theta}_j)^2$$

looks like a real mess but major simplification is imminent

Simplification: ①  $\hat{\theta}_j$  maximizes  $p(\theta_j | D, M_j, \mathcal{B})$  (15)

$$\text{so } \left. \left[ \frac{d}{d\theta_j} \log p(\theta_j | D, M_j, \mathcal{B}) \right]_{\theta_j = \hat{\theta}_j} = 0 \right\} \text{Assume } \textcircled{2} \text{ now that the prior}$$

is LI: then it will be close to flat

in the region around  $\theta_j = \hat{\theta}_j$ , so  $\left[ \frac{d}{d\theta_j} \log p(\theta_j | M_j, \mathcal{B}) \right]_{\theta_j = \hat{\theta}_j}$

$$\text{is also } \doteq 0$$

Therefore the entire <sup>linear</sup> term

$$\left[ \text{---} \right] (\theta_j - \hat{\theta}_j) \doteq 0$$

③ by the same reasoning as in ②,

$$\left[ \frac{d^2}{d\theta_j^2} \log p(\theta_j | D, M_j, \mathcal{B}) \right]_{\theta_j = \hat{\theta}_j} \text{ is also } \doteq 0$$

$$= -\mathbf{I}(\hat{\theta}_j)$$

$$\text{and } \left[ \frac{d^2}{d\theta_j^2} \log p(\theta_j | D, M_j, \mathcal{B}) \right]_{\theta_j = \hat{\theta}_j} = - \left( \text{observed information} \right) (!)$$

So Mr. Laplace's approximation becomes (16)  
(integral) =  $p(D | \theta_j; M_j; \mathcal{B}) \cdot p(\theta_j | n_j; \mathcal{B})$

$$= \exp \left[ \ell(\theta_j | D; M_j; \mathcal{B}) + \log p(\theta_j | n_j; \mathcal{B}) \right]$$

$$\cdot \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \hat{I}(\theta_j) (\theta_j - \hat{\theta}_j)^2 \right] \text{ here}$$

integrating over  $\theta_j$ , so the first two terms are constant in  $\theta_j$  & can be brought out in front of the integral

This yields  
(with  $k_j = 1$ )

$$\left( \text{evidence in favor of } M_j \right) = p(D | M_j; \mathcal{B}) \cdot$$

$$\int_{\mathcal{R}} p(D | \theta_j; M_j; \mathcal{B}) p(\theta_j | n_j; \mathcal{B}) d\theta_j$$

$$= \exp \left[ \ell(\hat{\theta}_j | D; M_j; \mathcal{B}) \right] \cdot \exp \left[ \log p(\hat{\theta}_j | n_j; \mathcal{B}) \right] \cdot (*)$$



in which  $\Theta = \int_{\mathcal{R}} \exp\left[-\frac{1}{2} \tilde{I}(\bar{\theta}_j) (\theta_j - \bar{\theta}_j)^2\right] d\theta_j$  (17)

But this is <sup>just</sup> the ~~integral of an~~ ~~un-normalized~~ Normal distribution (!) As Wolfram will tell you after some <sup>cooking</sup>

$$\int_{-\infty}^{\infty} \exp\left[-c_1 (\theta - c_2)^2\right] d\theta = \sqrt{\frac{\pi}{c_1}} \quad (c_1 > 0)$$

so

(evidence in favor of  $M_j$ ) =  $p(D | M_j; \mathcal{B})$

$$= \exp[\mathcal{L}(\bar{\theta}_j | D, M_j; \mathcal{B})] \cdot \exp[\log p(\bar{\theta}_j | M_j; \mathcal{B})] \cdot \sqrt{\frac{\pi}{\frac{1}{2} \tilde{I}(\bar{\theta}_j)}}$$

and therefore

•  $\log$  (evidence in favor of  $M_j$ ) =  $\log p(D | M_j; \mathcal{B})$

(\*\*)

$$= \boxed{\mathcal{L}(\bar{\theta}_j | D, M_j)} + \log p(\bar{\theta}_j | M_j; \mathcal{B}) + \frac{1}{2} \log(2\pi) - \frac{1}{2} \log |\tilde{I}(\bar{\theta}_j)|$$

an old friend

Looking back at the Taylor series approximation on p. (14),  $\log(\text{integrand}) =$

$\ell(\tilde{\theta}_j | D; \mathcal{B}) + \log p(\tilde{\theta}_j | M_j; \mathcal{B})$

$O(1)$  or  $n \uparrow$

+  ~~$\{ \text{complicated constant} \} (\theta - \tilde{\theta}_j)$~~

$O(\frac{1}{\sqrt{n}})$   
(think of  $SE(\tilde{\theta}_j) = \frac{c}{\sqrt{n}}$ )

+  $\{ \text{complicated constant} \} (\theta - \tilde{\theta}_j)^2$

$O(\frac{1}{n})$   
(think of  $\hat{V}(\tilde{\theta}_j)$ )

The generalization of (\*\*\*) on p. (17) turns out to be (including the  $O(\cdot)$  story):

~~Laplace approximation to log marginal likelihood~~

$\log(\text{evidence in favor of } M_j) = \log p(D | M_j; \mathcal{B}) = \ell(\tilde{\theta}_j | M_j; \mathcal{B})$

# of free parameters in  $M_j$

+  $\log p(\tilde{\theta}_j | M_j; \mathcal{B}) + \frac{k_j}{2} \log(2\pi) - \frac{1}{2} \log |\hat{I}(\tilde{\theta}_j)| + n$

$O(\frac{1}{n})$

in which  $|I(\hat{\theta}_j)|$  is the determinant of <sup>(19)</sup>  
the observed information matrix and  
 $\hat{\theta}_j$  is the ( $k_j$ -dimensional) MLE vector.

---

Since the <sup>Laplace</sup> approximation is accurate to  $O(\frac{1}{n})$ ,  
it's quite accurate; after all, the CLT  
is only accurate to  $O(\frac{1}{\sqrt{n}})$  (R code)

---

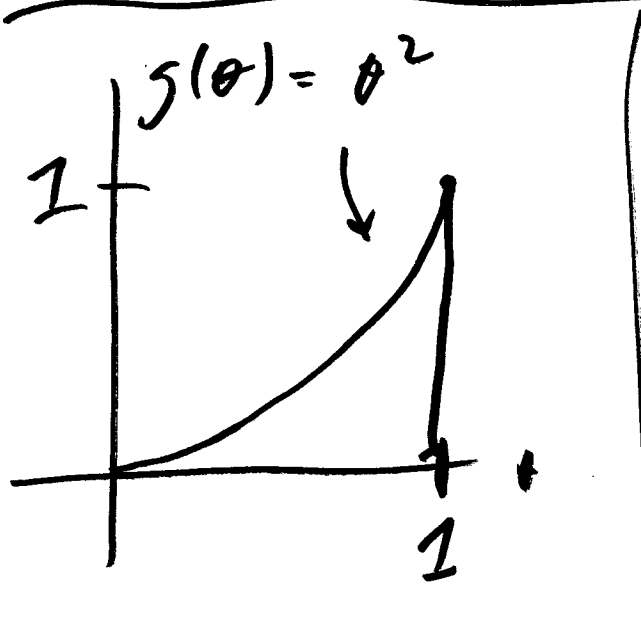
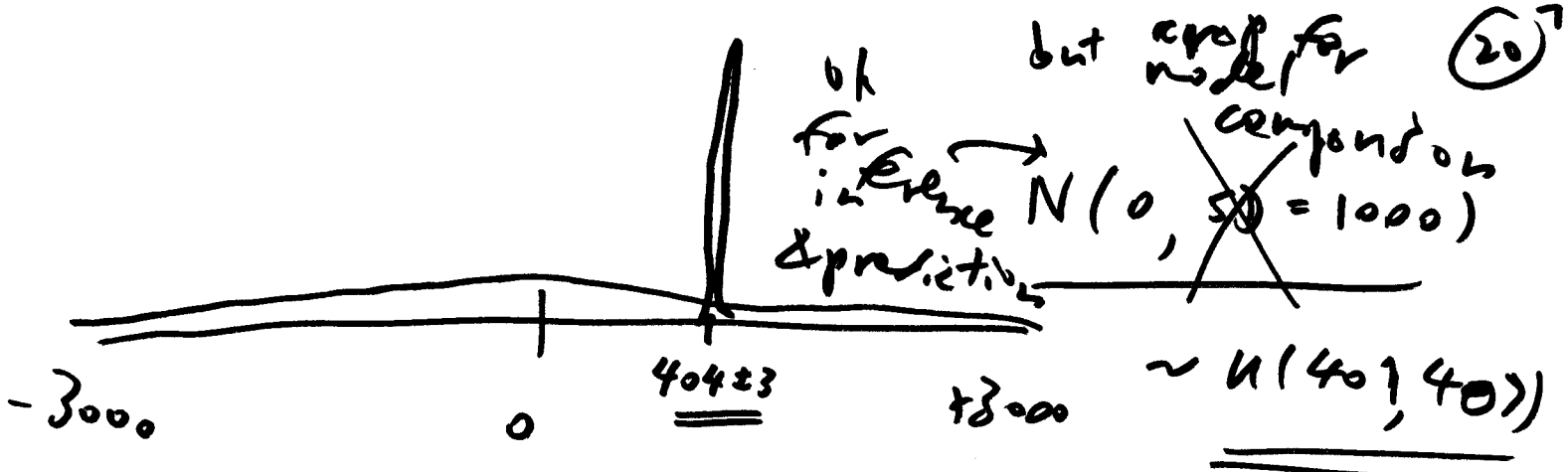
You can compute ~~\*\*\*~~ using only output  
from 'optim' applied to the log likelihood  
function with the 'hessian = T' option

---

One more wrinkle in this story:

---

One way to obtain useful LF priors for  
Bayesian model comparison with Bayes factors  
(to be continued)



$$\int_0^1 g(\theta) d\theta = \int_0^1 \theta^2 d\theta$$

$$= \frac{\theta^3}{3} \Big|_0^1 = \frac{1}{3}$$

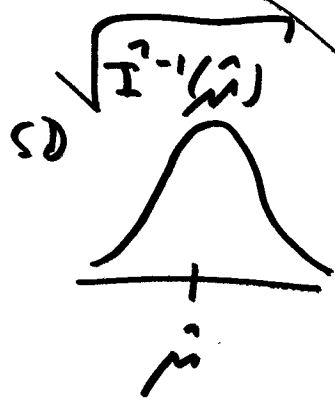
$\int_0^1 g(\theta) p(\theta) d\theta$

PJF

$\leftarrow U(0, 1)$

$\int_0^1 g(\theta) d\theta$

$$= \frac{1}{M} \sum_{i=1}^M g(\theta_i^*)$$



$p(\mu) M, \dots, B)$

n large

$\theta_i^* \sim U(0, 1)$

$$c_1 \exp[-c_2 (\mu - \hat{\mu})^2]$$