

Bayesian model selection and averaging

Will Penny

SPM for MEG/EEG, 15th May 2012

Bayes rule for
models

Bayes factors

Linear Models

Complexity

Nonlinear Models

Model Families

Model Averaging

Group Model
Inference

Fixed Effects

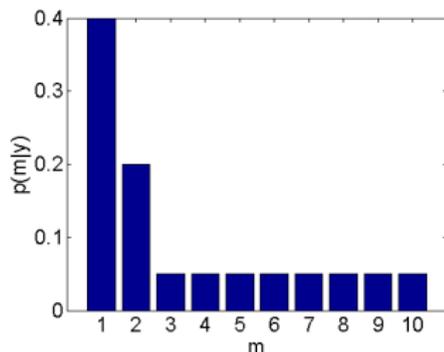
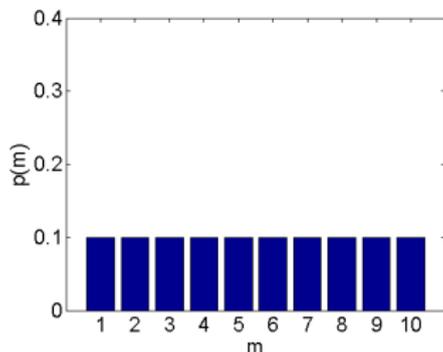
Random Effects

Gibbs Sampling

References

Bayes rule for models

A prior distribution over model space $p(m)$ (or 'hypothesis space') can be updated to a posterior distribution after observing data y .



This is implemented using Bayes rule

$$p(m|y) = \frac{p(y|m)p(m)}{p(y)}$$

where $p(y|m)$ is referred to as the evidence for model m and the denominator is given by

$$p(y) = \sum_{m'} p(y|m')p(m')$$

Bayes Factors

The Bayes factor for model j versus i is the ratio of model evidences

$$B_{ji} = \frac{p(y|m=j)}{p(y|m=i)}$$

We have

$$B_{ij} = \frac{1}{B_{ji}}$$

Posterior Model Probability

Given equal priors, $p(m = i) = p(m = j)$ the posterior model probability is

$$\begin{aligned} p(m = i|y) &= \frac{p(y|m = i)}{p(y|m = i) + p(y|m = j)} \\ &= \frac{1}{1 + \frac{p(y|m=j)}{p(y|m=i)}} \\ &= \frac{1}{1 + B_{ji}} \\ &= \frac{1}{1 + \exp(\log B_{ji})} \\ &= \frac{1}{1 + \exp(-\log B_{ij})} \end{aligned}$$

Posterior Model Probability

Hence

$$p(m = i|y) = \sigma(\log B_{ij})$$

where B_{ij} is the Bayes factor for model i versus model j and

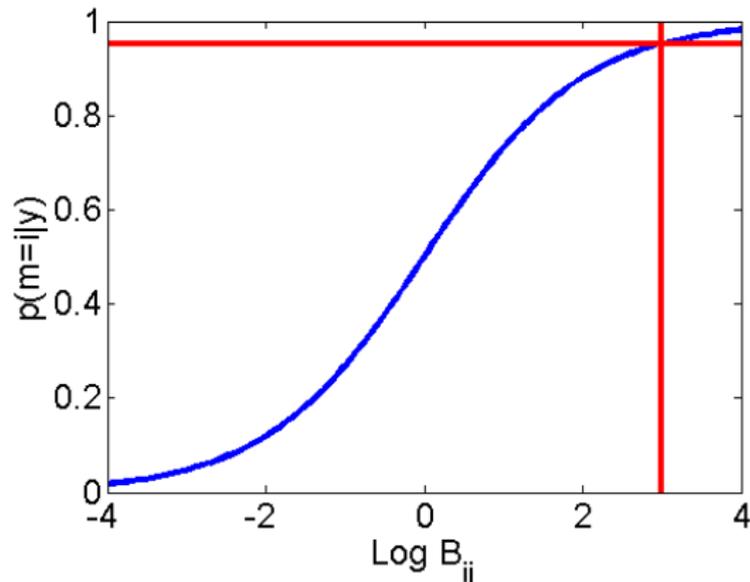
$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

is the sigmoid function.

Bayes factors

The posterior model probability is a sigmoidal function of the log Bayes factor

$$p(m = i | y) = \sigma(\log B_{ij})$$



Bayes factors

The posterior model probability is a sigmoidal function of the log Bayes factor

$$p(m = i|y) = \sigma(\log B_{ij})$$

Table 1
Interpretation of Bayes factors

B_{ij}	$p(m = i y)$ (%)	Evidence in favor of model i
1–3	50–75	Weak
3–20	75–95	Positive
20–150	95–99	Strong
≥ 150	≥ 99	Very strong

Bayes factors can be interpreted as follows. Given candidate hypotheses i and j , a Bayes factor of 20 corresponds to a belief of 95% in the statement ‘hypothesis i is true’. This corresponds to strong evidence in favor of i .

From Raftery (1995).

Odds Ratios

If we don't have uniform priors one can work with odds ratios.

The prior and posterior odds ratios are defined as

$$\pi_{ij}^0 = \frac{p(m=i)}{p(m=j)}$$
$$\pi_{ij} = \frac{p(m=i|y)}{p(m=j|y)}$$

respectively, and are related by the Bayes Factor

$$\pi_{ij} = B_{ij} \times \pi_{ij}^0$$

eg. priors odds of 2 and Bayes factor of 10 leads posterior odds of 20.

An odds ratio of 20 is 20-1 ON in bookmakers parlance.

Model Evidence

The model evidence is not, in general, straightforward to compute since computing it involves integrating out the dependence on model parameters

$$\begin{aligned} p(y|m) &= \int p(y, \theta|m) d\theta \\ &= \int p(y|\theta, m) p(\theta|m) d\theta \end{aligned}$$

Because we have marginalised over θ the evidence is also known as the marginal likelihood.

But for linear, Gaussian models there is an analytic solution.

Linear Models

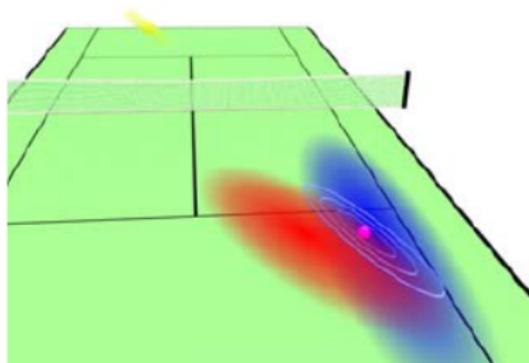
For Linear Models

$$y = Xw + e$$

where X is a design matrix and w are now regression coefficients. For prior mean μ_w , prior covariance C_w , observation noise covariance C_y the posterior distribution is given by

$$S_w^{-1} = X^T C_y^{-1} X + C_w^{-1}$$

$$m_w = S_w \left(X^T C_y^{-1} y + C_w^{-1} \mu_w \right)$$



Bayesian model selection and averaging

Will Penny

Bayes rule for models

Bayes factors

Linear Models

Complexity

Nonlinear Models

Model Families

Model Averaging

Group Model Inference

Fixed Effects

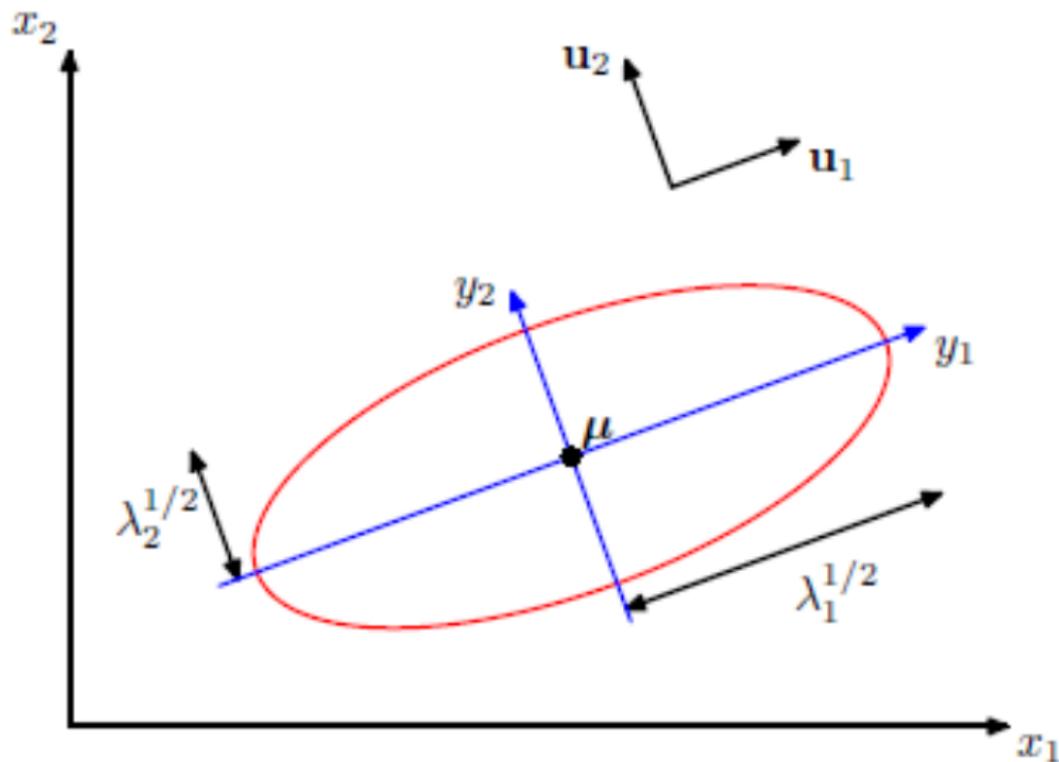
Random Effects

Gibbs Sampling

References

Covariance matrices

The determinant of a covariance matrix, $|C|$, measures the volume.



Model Evidence

The log model evidence comprises sum squared precision weighted prediction errors and Occam factors

$$L = -\frac{1}{2} e_y^T C_y^{-1} e_y - \frac{1}{2} \log |C_y| - \frac{N_y}{2} \log 2\pi \\ - \frac{1}{2} e_w^T C_w^{-1} e_w - \frac{1}{2} \log \frac{|C_w|}{|S_w|}$$

where prediction errors are the difference between what is expected and what is observed

$$e_y = y - X m_w$$

$$e_w = m_w - \mu_w$$

See Bishop (2006) for derivation.

Accuracy and Complexity

The log evidence for model m can be split into an accuracy and a complexity term

$$L(m) = \text{Accuracy}(m) - \text{Complexity}(m)$$

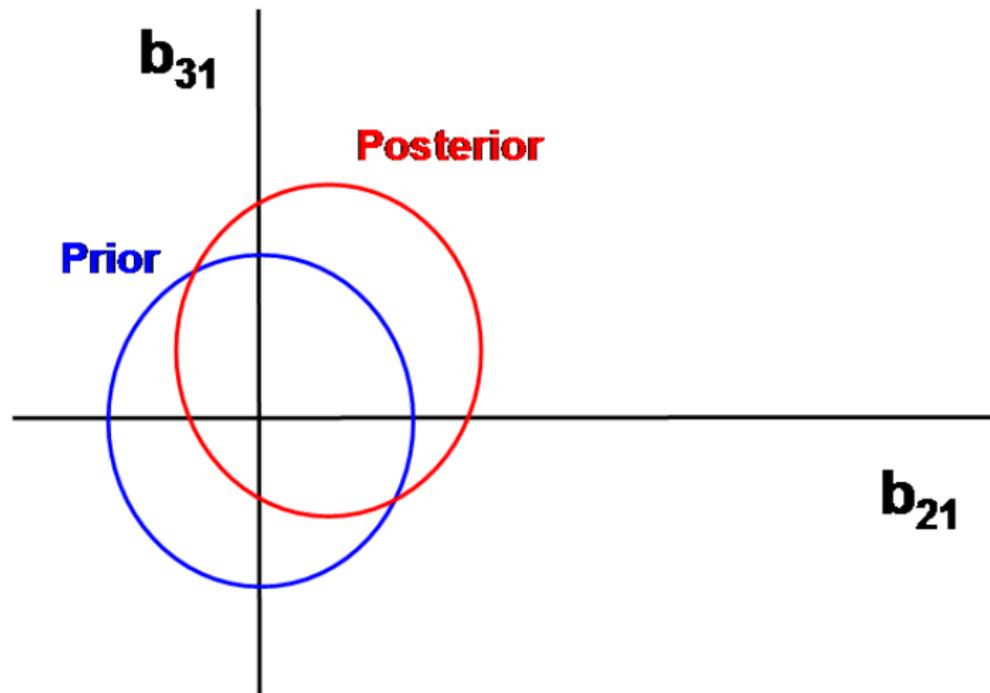
where

$$\text{Accuracy}(m) = -\frac{1}{2} \mathbf{e}_y^T \mathbf{C}_y^{-1} \mathbf{e}_y - \frac{1}{2} \log |\mathbf{C}_y| - \frac{N_y}{2} \log 2\pi$$

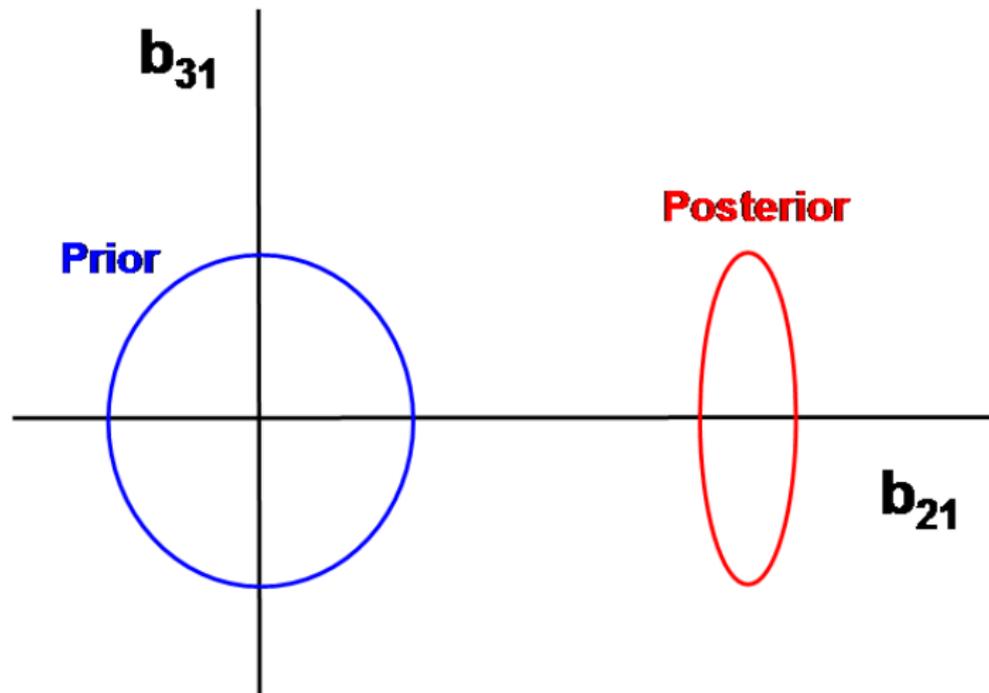
and

$$\begin{aligned} \text{Complexity}(m) &= \frac{1}{2} \mathbf{e}_w^T \mathbf{C}_w^{-1} \mathbf{e}_w + \frac{1}{2} \log \frac{|\mathbf{C}_w|}{|\mathbf{S}_w|} \\ &\approx \text{KL}(\text{prior} || \text{posterior}) \end{aligned}$$

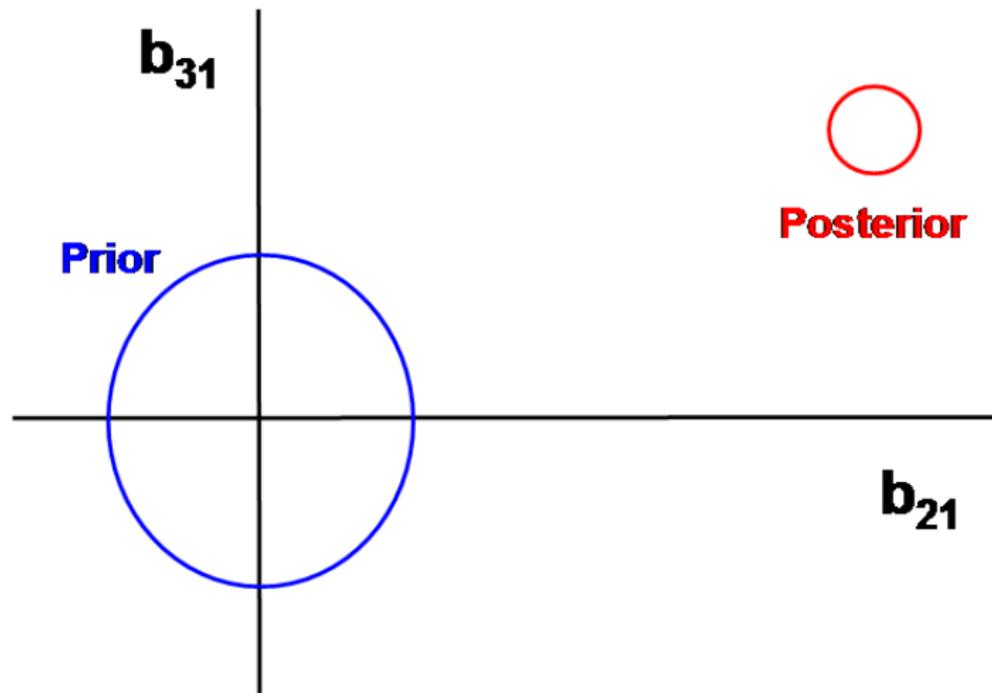
Small KL



Medium KL



Big KL



Nonlinear Models

For nonlinear models, we replace the true posterior with the approximate posterior (m_w, S_w) , and the previous expression becomes an approximation to the log model evidence called the (negative) Free Energy

$$F = -\frac{1}{2} e_y^T C_y^{-1} e_y - \frac{1}{2} \log |C_y| - \frac{N_y}{2} \log 2\pi \\ - \frac{1}{2} e_w^T C_w^{-1} e_w - \frac{1}{2} \log \frac{|C_w|}{|S_w|}$$

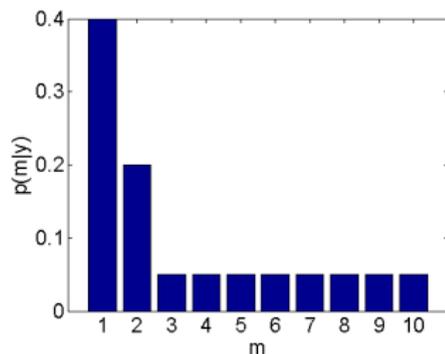
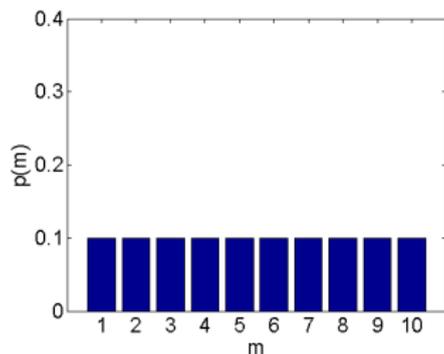
where

$$e_y = y - g(m_w) \\ e_w = m_w - \mu_w$$

and $g(m_w)$ is the DCM prediction. This is used to approximate the model evidence for DCMs (see Penny, Neuroimage, 2011 for more).

Bayes rule for models

A prior distribution over model space $p(m)$ (or 'hypothesis space') can be updated to a posterior distribution after observing data y .

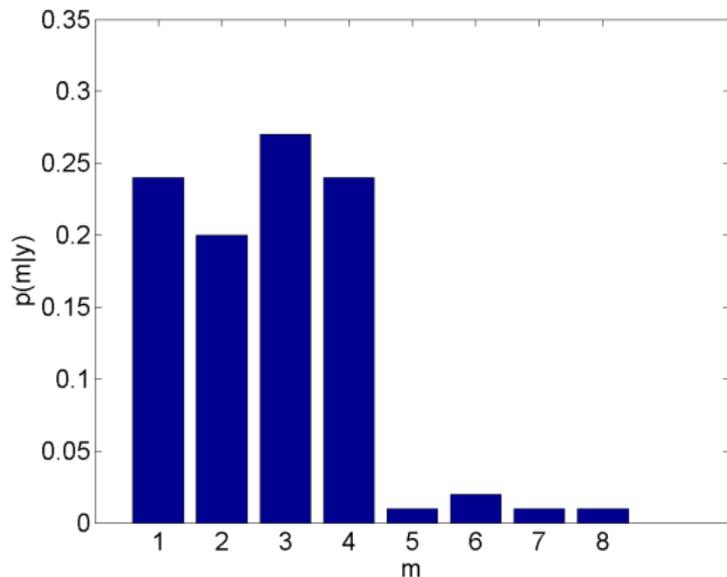


This is implemented using Bayes rule

$$p(m|y) = \frac{p(y|m)p(m)}{p(y)}$$

Posterior Model Probabilities

Say we've fitted 8 DCMs and get the following distribution over models

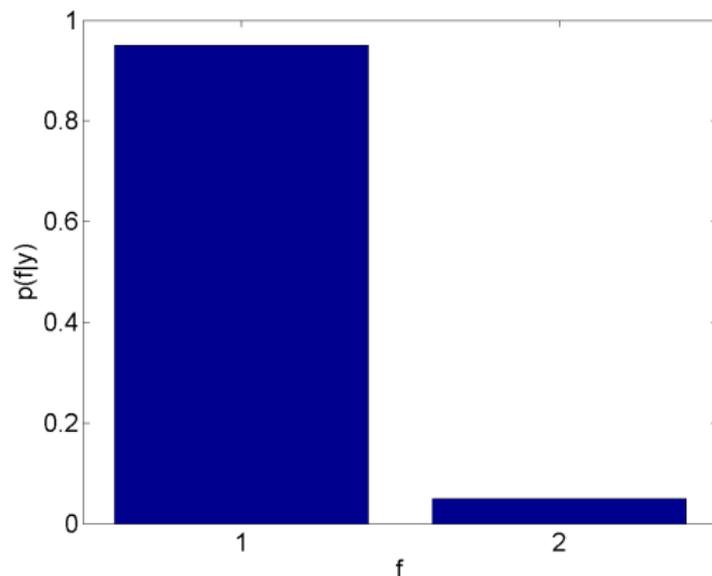


Similar models share probability mass (dilution). The probability for any single model can become very small esp. for large model spaces.

Model Families

Assign model m to family f eg. first four to family one, second four to family two. The posterior family probability is then

$$p(f|y) = \sum_{m \in S_f} p(m|y)$$



Different Sized Families

If we have K families, then to avoid bias in family inference we wish to have a uniform prior at the family level

$$p(f) = \frac{1}{K}$$

The prior family probability is related to the prior model probability

$$p(f) = \sum_{m \in S_f} p(m)$$

where the sum is over all N_f models in family f . So we set

$$p(m) = \frac{1}{KN_f}$$

for all models in family f before computing $p(m|y)$. This allows us to have families with unequal numbers of models.

Different Sized Families

So say we have two families. We want a prior for each family of $p(f) = 0.5$.

If family one has $N_1 = 2$ models and family two has $N_2 = 8$ models, then we set

$$p(m) = \frac{1}{2} \times \frac{1}{2} = 0.25$$

for all models in family one and

$$p(m) = \frac{1}{2} \times \frac{1}{8} = 0.0625$$

for all models in family two.

These are then used in Bayes rule for models

$$p(m|y) = \frac{p(y|m)p(m)}{p(y)}$$

Model Averaging

Each DCM.mat file stores the posterior mean (DCM.Ep) and covariance (DCM.Cp) for each fitted model. This defines the posterior mean over parameters for that model, $p(\theta|m, y)$.

This can then be combined with the posterior model probabilities $p(m|y)$ to compute a posterior over parameters

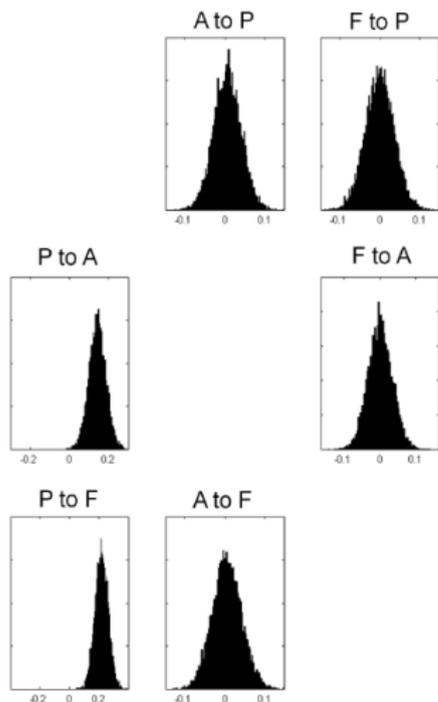
$$\begin{aligned} p(\theta|y) &= \sum_m p(\theta, m|y) \\ &= \sum_m p(\theta|m, y)p(m|y) \end{aligned}$$

which is independent of model assumptions (within the chosen set). Here, we marginalise over m .

The sum over m could be restricted to eg. models within the winning family.

Model Averaging

The distribution $p(\theta|y)$ can be gotten by sampling;
sample m from $p(m|y)$, then sample θ from $p(\theta|m, y)$.



If a connection doesn't exist for model m the relevant samples are set to zero.

Group Parameter Inference

Bayesian model
selection and
averaging

Will Penny

If i th subject has posterior mean value m_i we can use these in Summary Statistic approach for group parameter inference (eg two-sample t-tests for control versus patient inferences).

eg P to A connection in controls: 0.20, 0.12, 0.32, 0.11, 0.01, ...

eg P to A connection in patients: 0.50, 0.42, 0.22, 0.71, 0.31, ...

Two sample t-test shows the P to A connection is stronger in patients than controls ($p < 0.05$).

Bayes rule for
models

Bayes factors

Linear Models

Complexity

Nonlinear Models

Model Families

Model Averaging

Group Model
Inference

Fixed Effects

Random Effects

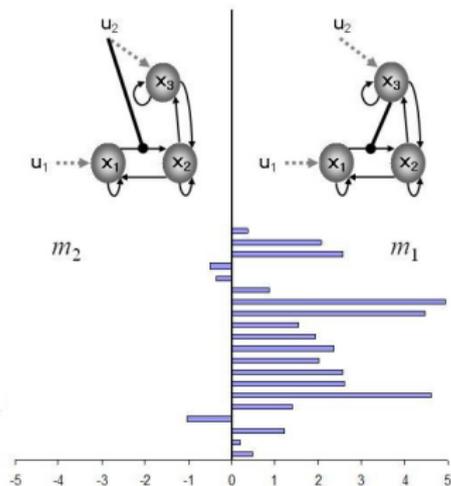
Gibbs Sampling

References

Fixed Effects

Two models, twenty subjects.

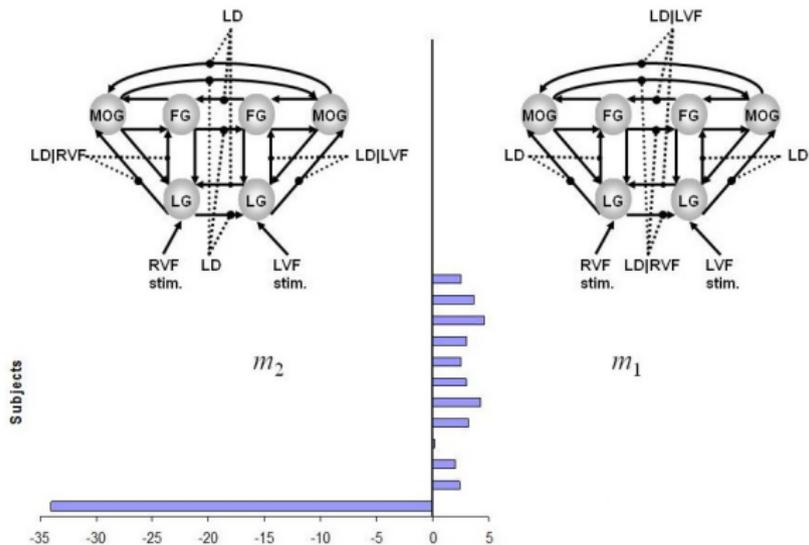
$$\log p(Y|m) = \sum_{n=1}^N \log p(y_n|m)$$



The Group Bayes Factor (GBF) is

$$B_{ij} = \prod_{n=1}^N B_{ij}(n)$$

Random Effects



11/12=92% subjects favour model 1.

$GBF = 15$ in favour of model 2. FFX inference does not agree with the majority of subjects.

Random Effects

For RFX analysis it is possible that different subjects use different models. If we knew exactly which subjects used which models then this information could be represented in a $[N \times M]$ assignment matrix, A , with entries $a_{nm} = 1$ if subject m used model n , and $a_{nm} = 0$ otherwise.

For example, the following assignment matrix

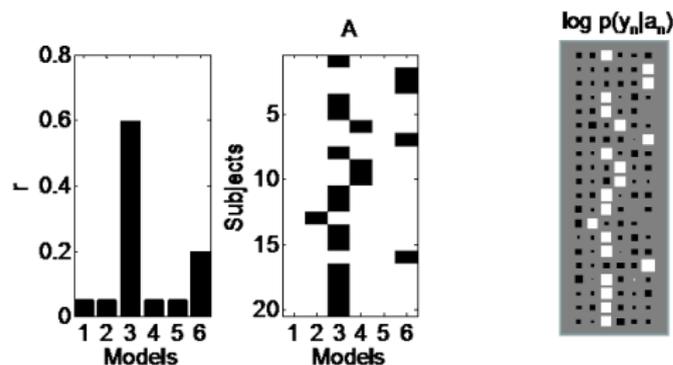
$$A = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

indicates that subjects 1 and 2 used model 2 and subject 3 used model 1.

We denote r_m as the frequency with which model m is used in the population. We also refer to r_m as the model probability.

Generative Model

In our generative model we have a prior $p(r|\alpha)$. A vector of probabilities is then drawn from this.



An assignment for each subject a_n is then drawn from $p(a_n|r)$. Finally a_n specifies which log evidence value to use for each subject. This specifies $p(y_n|a_n)$.

The joint likelihood for the RFX model is

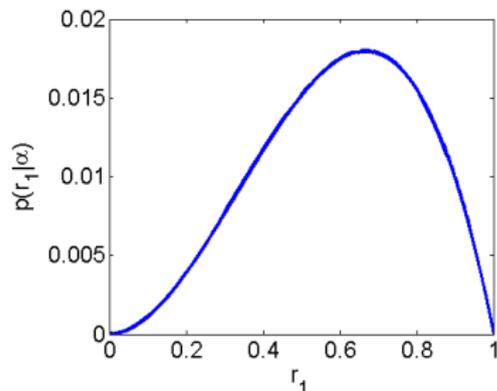
$$p(y, a, r|\alpha) = \prod_{n=1}^N [p(y_n|a_n)p(a_n|r)]p(r|\alpha)$$

Prior Model Frequencies

We define a prior distribution over r which is a Dirichlet

$$p(r|\alpha_0) = \text{Dir}(\alpha_0) = \frac{1}{Z} \prod_{m=1}^M r_m^{\alpha_0(m)-1}$$

where Z is a normalisation term and the parameters, α_0 , are strictly positively valued and the m th entry can be interpreted as the number of times model m has been selected.



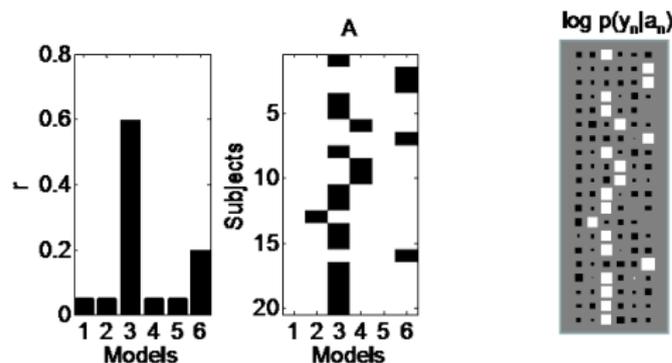
Example with $\alpha_0 = [3, 2]$ and $r = [r_1, 1 - r_1]$.

In the RFX generative model we use a uniform prior $\alpha_0 = [1, 1]$ or more generally $\alpha_0 = \text{ones}(1, M)$.

Model Assignment

The probability of the 'assignment vector', a_n , is then given by the multinomial density

$$p(a_n|r) = \text{Mult}(r) = \prod_{m=1}^M r_m^{a_{nm}}$$



The assignments then indicate which entry in the model evidence table to use for each subject, $p(y_n|a_n)$.

Gibbs Sampling

Samples from the posterior densities $p(r|y)$ and $p(a|y)$ can be drawn using Gibbs sampling (Gelman et al 1995).

This can be implemented by alternately sampling from

$$r \sim p(r|a, y)$$

$$a \sim p(a|r, y)$$

and discarding samples before convergence.

This is like a sample-based EM algorithm.

Gibbs Sampling

STEP 1: model probabilities are drawn from the prior distribution

$$r \sim \text{Dir}(\alpha_{\text{prior}})$$

where by default we set $\alpha_{\text{prior}}(m) = \alpha_0$ for all m (but see later).

STEP 2: For each subject $n = 1..N$ and model $m = 1..M$ we use the model evidences from model inversion to compute

$$u_{nm} = \exp(\log p(y_n|m) + \log r_m)$$

$$g_{nm} = \frac{u_{nm}}{\sum_{m=1}^M u_{nm}}$$

Here, g_{nm} is our posterior belief that model m generated the data from subject n .

Gibbs Sampling

STEP 3: For each subject, model assignment vectors are then drawn from the multinomial distribution

$$\mathbf{a}_n \sim \text{Mult}(\mathbf{g}_n)$$

We then compute new model counts

$$\beta_m = \sum_{n=1}^N \mathbf{a}_{nm}$$
$$\alpha_m = \alpha_{\text{prior}}(m) + \beta_m$$

and draw new model probabilities

$$r \sim \text{Dir}(\alpha)$$

Go back to STEP 2 !

Gibbs Sampling

Bayesian model
selection and
averaging

Will Penny

Bayes rule for
models

Bayes factors

Linear Models

Complexity

Nonlinear Models

Model Families

Model Averaging

Group Model
Inference

Fixed Effects

Random Effects

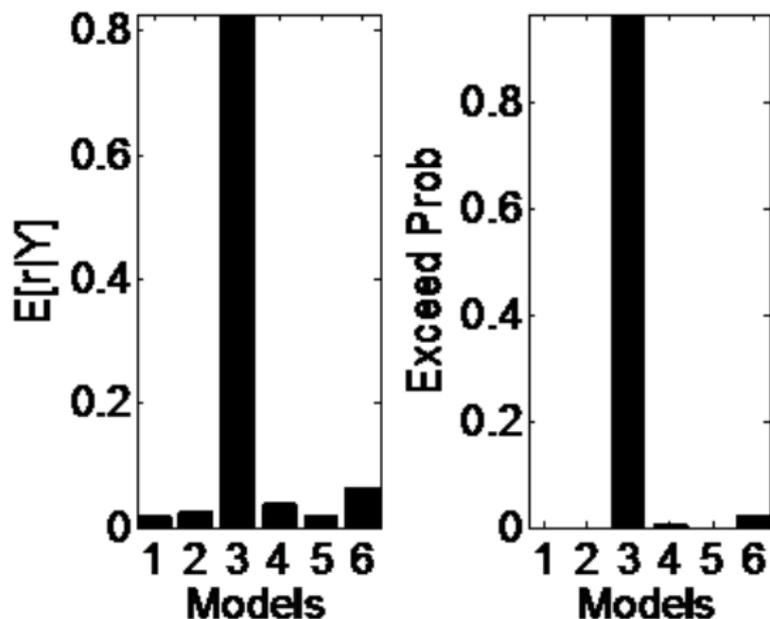
Gibbs Sampling

References

These steps are repeated N_d times. For the following results we used a total of $N_d = 20,000$ samples and discarded the first 10,000.

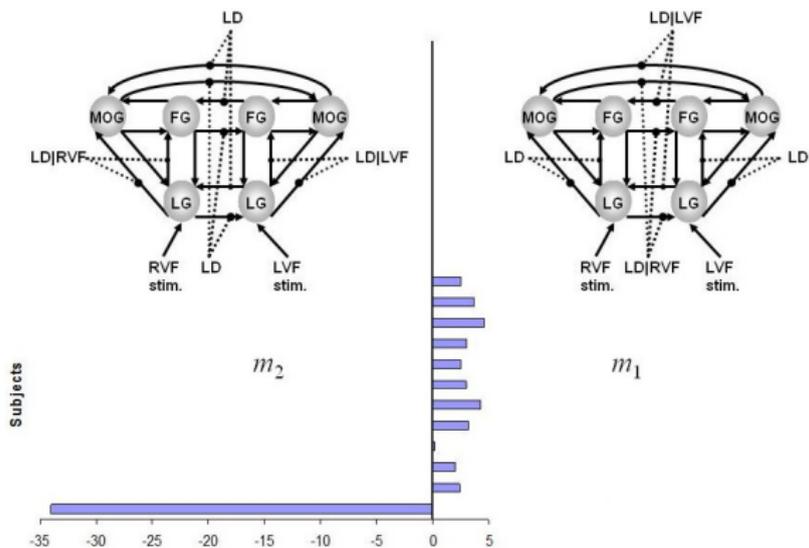
Gibbs Sampling

These remaining samples then constitute our approximation to the posterior distribution $p(r|Y)$. From this density we can compute usual quantities such as the posterior expectation, $E[r|Y]$.



Random Effects

11/12=92% subjects favoured model 1.



$$E[r_1 | Y] = 0.84$$

$$p(r_1 > r_2 | Y) = 0.99$$

where the latter is called the exceedance probability.

Dependence on Comparison Set

The ranking of models from RFX inference can depend on the comparison set.

Say we have two models with 7 subjects preferring model 1 and 10 ten subjects preferring model 2. The model frequencies are $r_1 = 7/17 = 0.41$ and $r_2 = 10/17 = 0.59$.

Now say we add a third model which is similar to the second, and that 4 of the subjects that used to prefer model 2 now prefer model 3. The model frequencies are now $r_1 = 7/17 = 0.41$, $r_2 = 6/17 = 0.35$ and $r_3 = 4/17 = 0.24$.

This is like voting in elections.

References

C. Bishop (2006) Pattern Recognition and Machine Learning. Springer.

A. Gelman et al. (1995) Bayesian Data Analysis. Chapman and Hall.

W. Penny (2011) Comparing Dynamic Causal Models using AIC, BIC and Free Energy. Neuroimage Available online 27 July 2011.

W. Penny et al (2010) Comparing Families of Dynamic Causal Models. PLoS CB, 6(3).

A Raftery (1995) Bayesian model selection in social research. In Marsden, P (Ed) Sociological Methodology, 111-196, Cambridge.

K Stephan et al (2009). Bayesian model selection for group studies. Neuroimage, 46(4):1004-17

Bayesian model selection and averaging

Will Penny

Bayes rule for models

Bayes factors

Linear Models

Complexity

Nonlinear Models

Model Families

Model Averaging

Group Model Inference

Fixed Effects

Random Effects

Gibbs Sampling

References