

Bayesian inference and Bayesian model selection

Klaas Enno Stephan



Translational Neuromodeling Unit



Universität
Zürich^{UZH}

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Lecture as part of "Methods & Models for fMRI data analysis",
University of Zurich & ETH Zurich, 26 November 2019

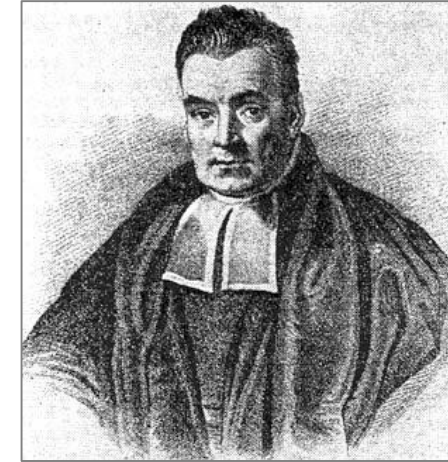
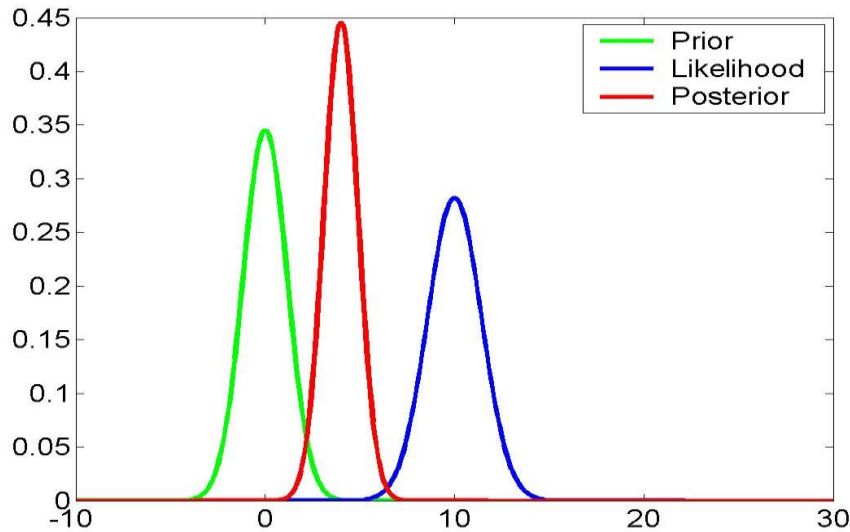
With slides from and many thanks to:

Kay Brodersen,

Will Penny,

Sudhir Shankar Raman

Bayes' rule



The Reverend Thomas Bayes
(1702-1761)

Likelihood \times prior: generative model

$$p(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y})}$$

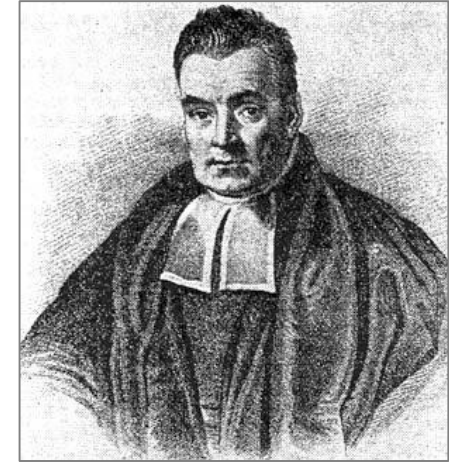
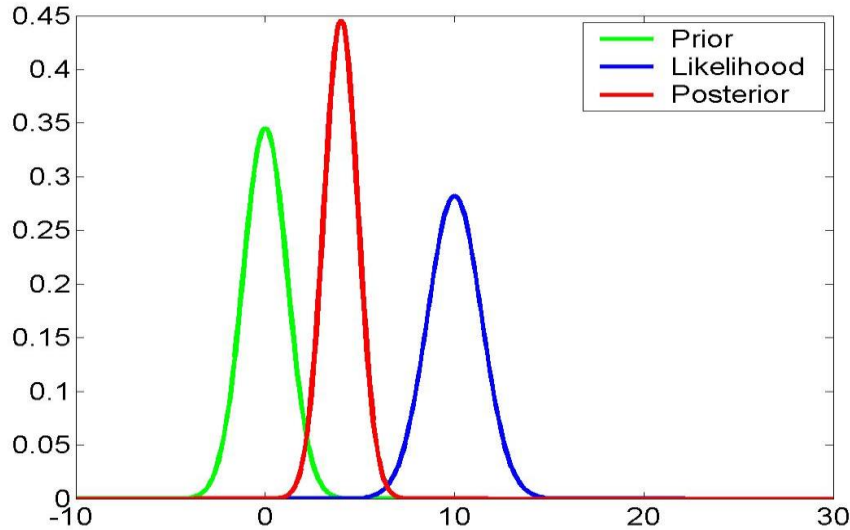
$\boldsymbol{\theta}$: parameters
 \mathbf{y} : data

Model evidence: normalisation
term and index for model goodness

"... the theorem expresses how a ... degree of belief should rationally change to account for availability of related evidence."

Wikipedia

Bayes' rule

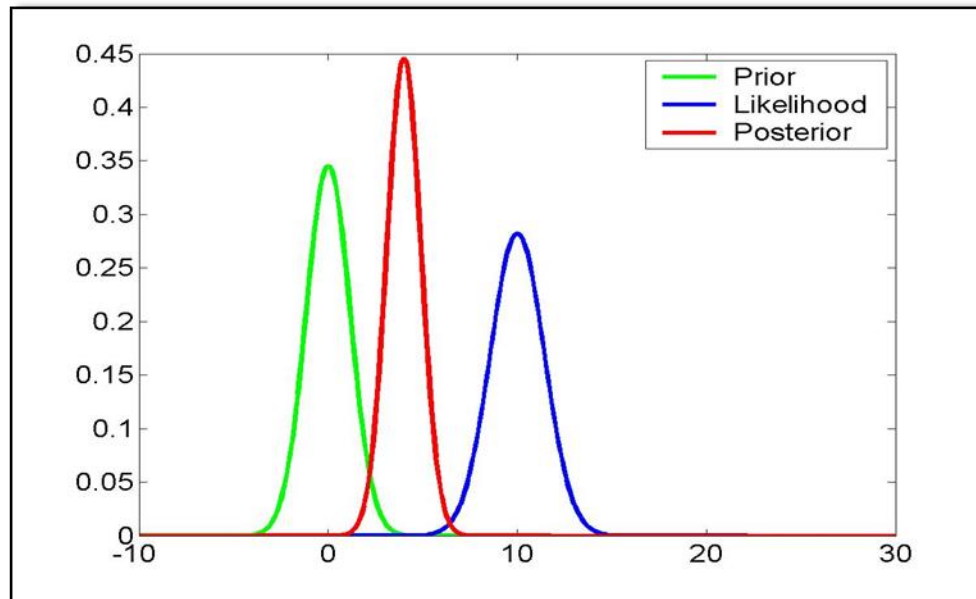


The Reverend Thomas Bayes
(1702-1761)

$$p(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y})}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Bayesian inference: an animation



Code courtesy by Guillaume Flandin

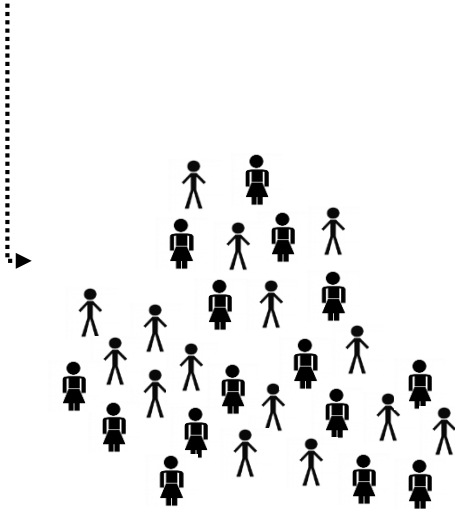
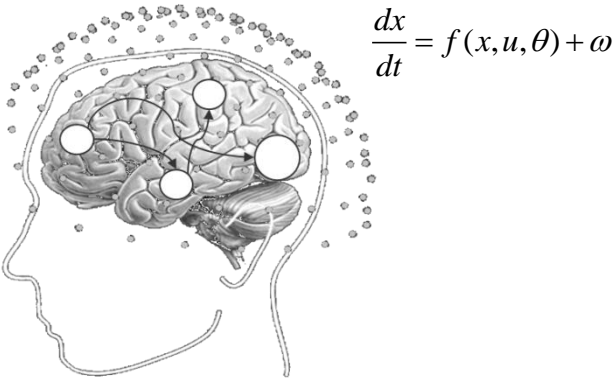
But why should I learn about Bayesian inference?

Because Bayesian principles are fundamental for

- **statistical inference** in general
- **system identification**
- **translational neuromodeling** ("computational assays")
 - computational psychiatry
 - computational neurology
 - computational psychosomatics
- contemporary **theories of brain function** (the "Bayesian brain")
 - predictive coding
 - free energy principle
 - active inference

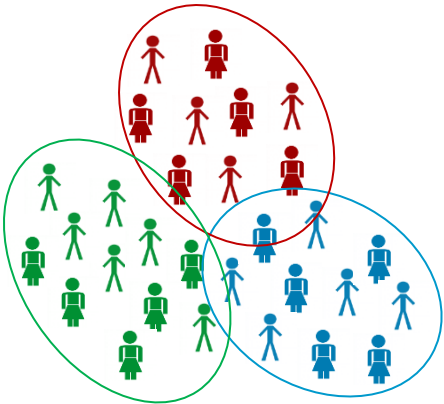
Translational Neuromodeling

1 Computational assays: Models of disease mechanisms



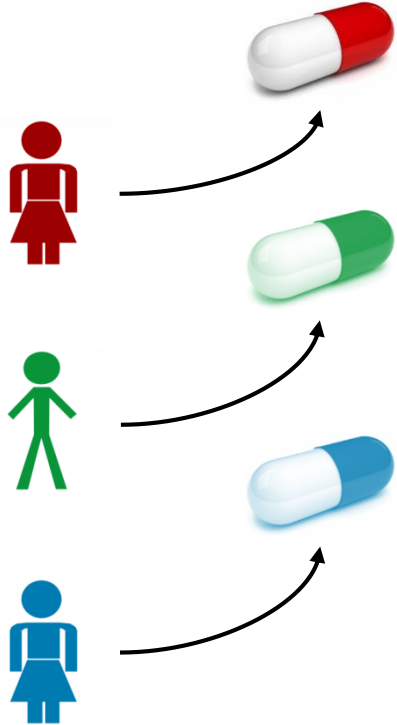
2 Application to brain activity and behaviour of individual patients

3 Detecting physiological subgroups (based on inferred mechanisms)

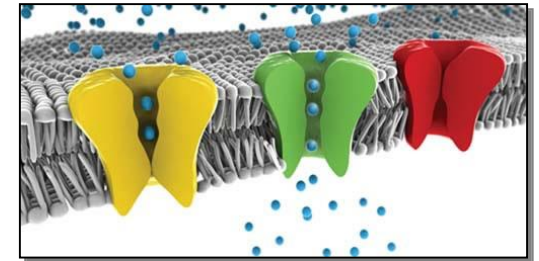
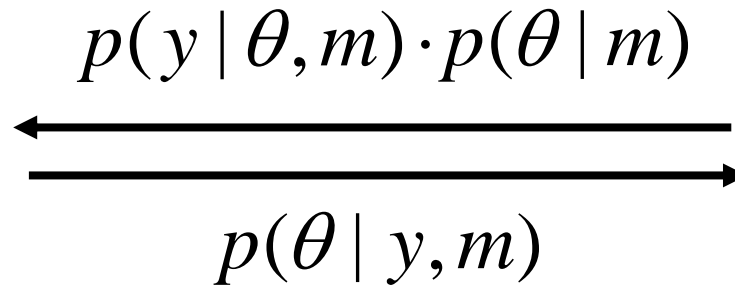
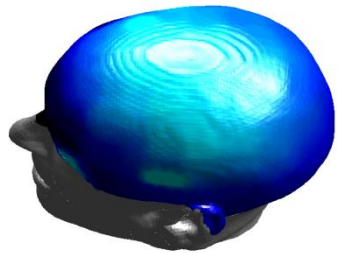
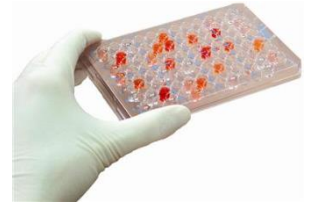


- disease mechanism A
- disease mechanism B
- disease mechanism C

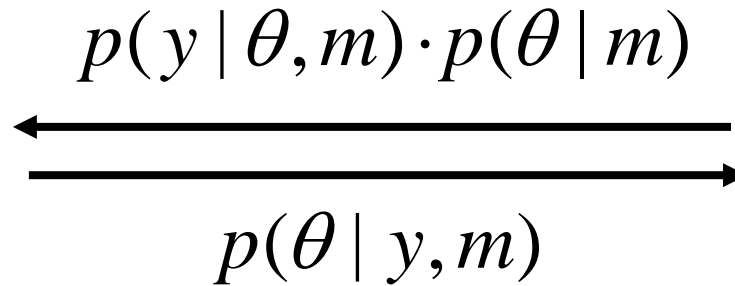
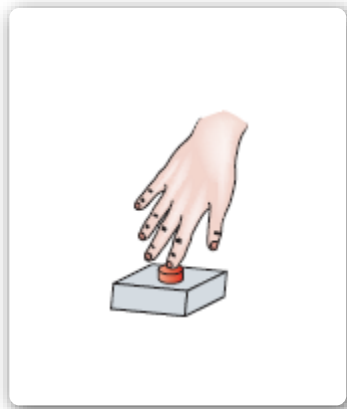
4 Individual treatment prediction



Generative models as "computational assays"



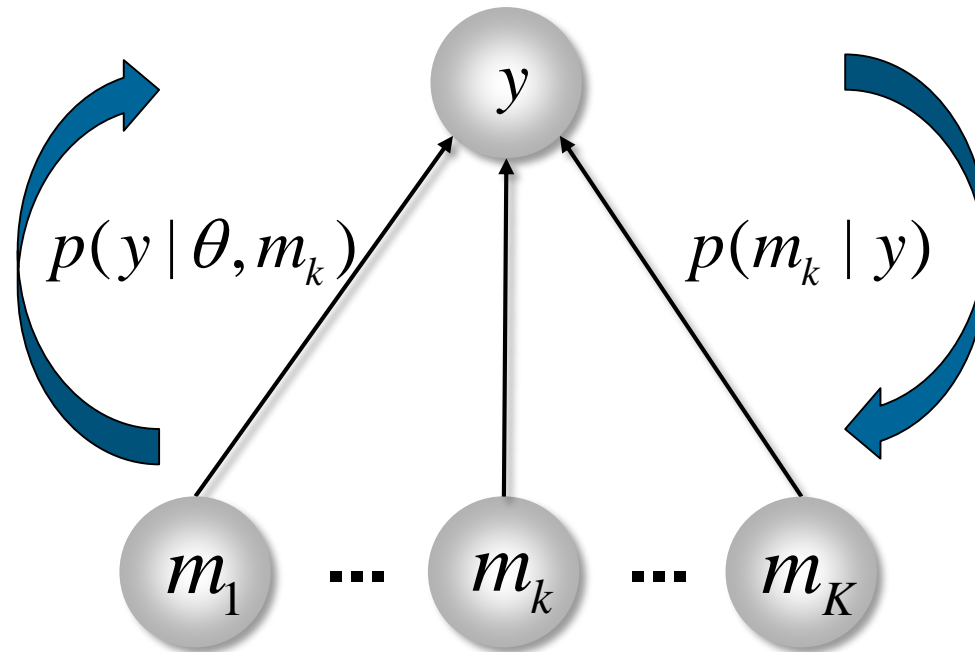
y = data, θ = parameters, m = model



Differential diagnosis by model selection

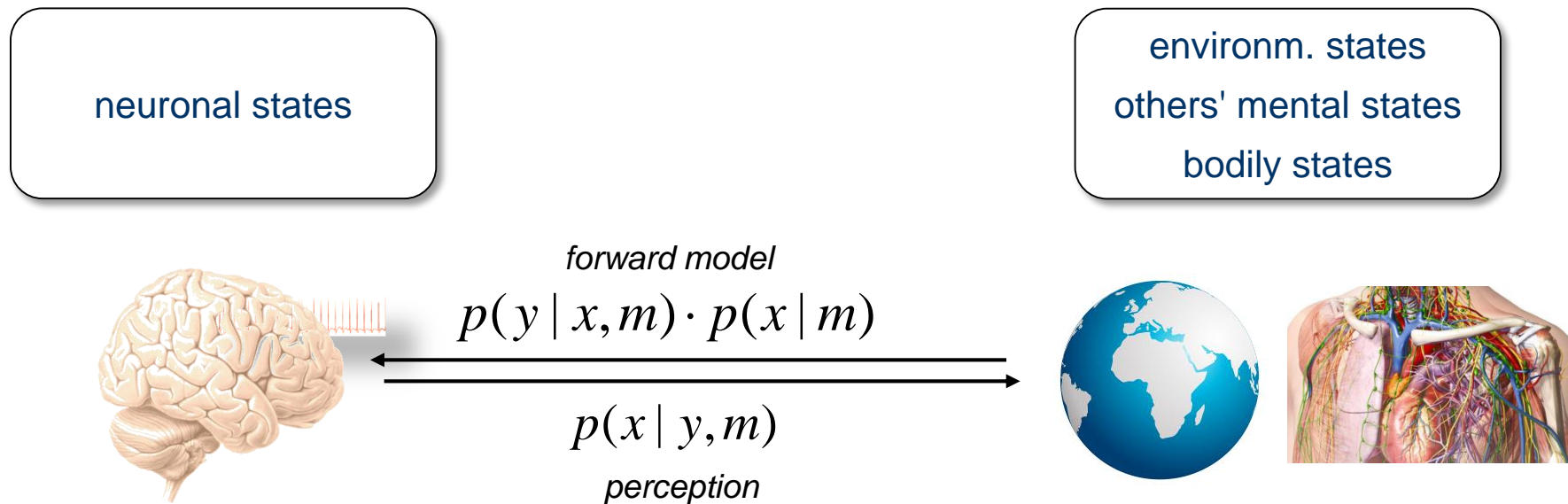
SYMPTOM
(behaviour
or physiology)

**HYPOTHETICAL
MECHANISM**



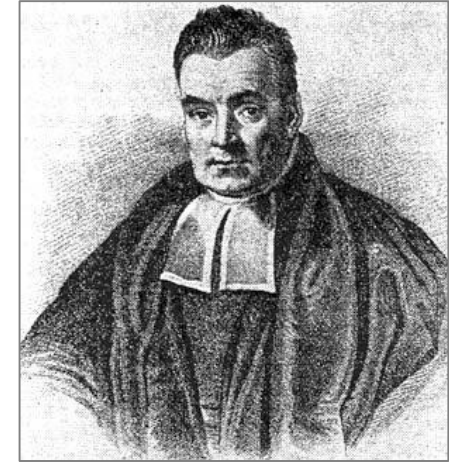
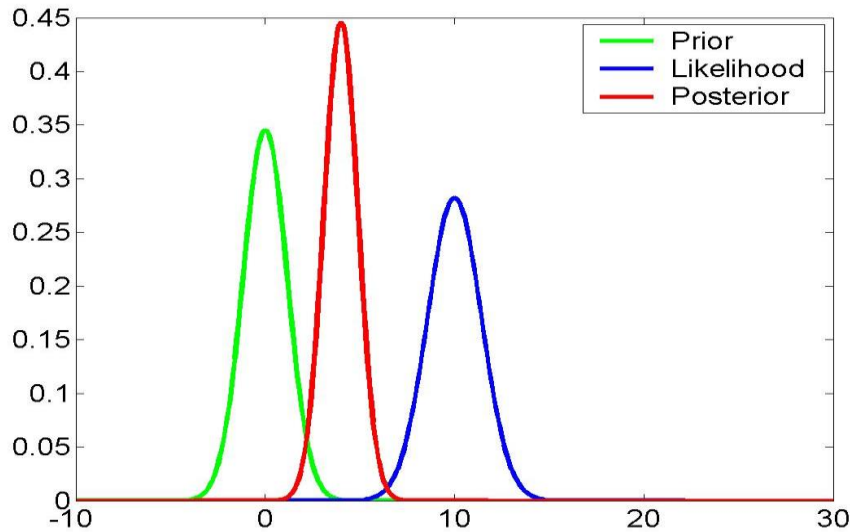
$$p(m_k | y) = \frac{p(y | m_k) p(m_k)}{\sum_k p(y | m_k) p(m_k)}$$

Perception = inversion of a hierarchical generative model



Back to the technicalities...

Bayes' rule

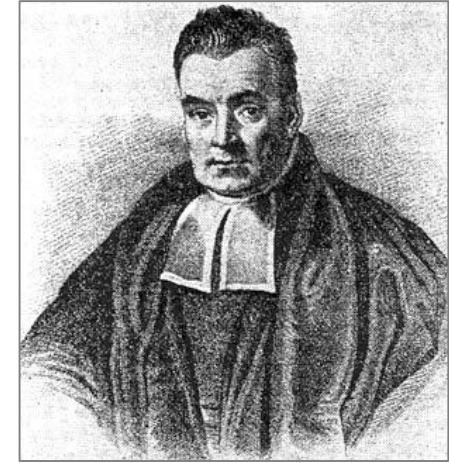
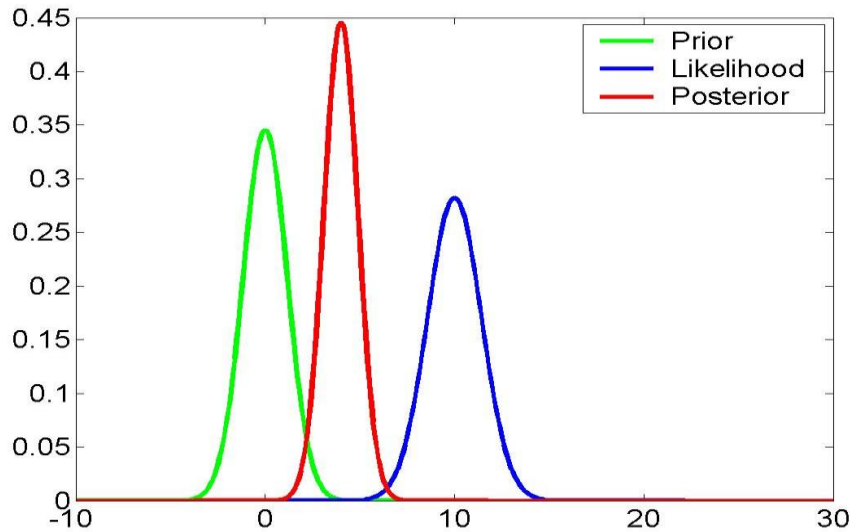


The Reverend Thomas Bayes
(1702-1761)

$$p(\boldsymbol{\theta} | \mathbf{y}, m) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, m) p(\boldsymbol{\theta} | m)}{p(\mathbf{y} | m)}$$

No change – just making the choice of a particular model m explicit.

Bayes' rule



The Reverend Thomas Bayes
(1702-1761)

$$p(\boldsymbol{\theta} | \mathbf{y}, m) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, m) p(\boldsymbol{\theta} | m)}{\int p(\mathbf{y} | \boldsymbol{\theta}, m) p(\boldsymbol{\theta} | m)}$$

posterior = likelihood • prior / evidence

Evidence:

probability that data were generated by model m , averaging over all possible parameter values (as weighted by the prior).

The evidence term

continuous θ

$$p(\theta | y) = \frac{p(y | \theta) p(\theta)}{\int p(y | \theta) p(\theta)}$$

discrete θ

$$p(\theta | y) = \frac{p(y | \theta) p(\theta)}{\sum_{\theta \in \Theta} p(y | \theta) p(\theta)}$$

Bayesian inference: A clinical example

- *"The probability of breast cancer is 1% for women aged forty who participate in routine screening. If a woman has breast cancer, the probability is 80% that she will get a positive mammogram. If a woman does not have breast cancer, the probability is 9.6% that she will also get a positive mammogram. A woman in this age group has a positive mammogram in a routine screening. What is the probability that she actually has breast cancer?"* (Gigerenzer & Hoffrage 1995)
- From this information, we can deduce:
 - $p(C+) = 0.01 \rightarrow p(C-) = 0.99$
 - $p(M+|C+) = 0.8$
 - $p(M+|C-) = 0.096$
- We can now apply Bayes' rule to compute the posterior probability:

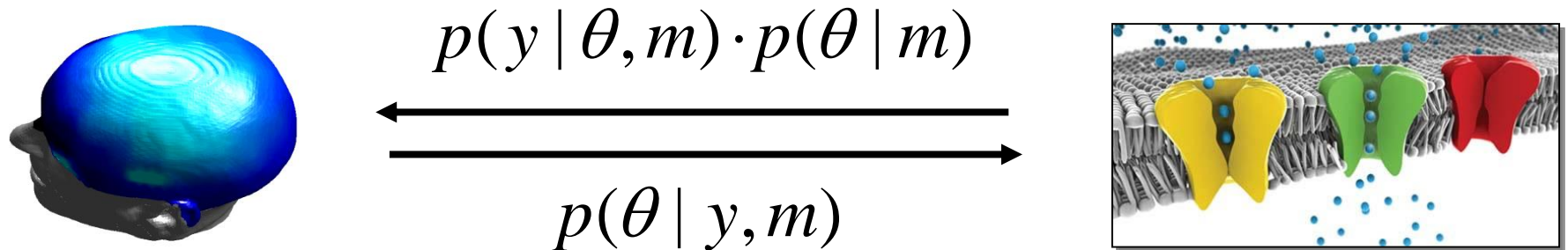
$$p(C+ | M+) = \frac{p(M+ | C+)p(C+)}{p(M+ | C+)p(C+) + p(M+ | C-)p(C-)}$$

Bayesian inference: A clinical example

- *"The probability of breast cancer is 1% for women aged forty who participate in routine screening. If a woman has breast cancer, the probability is 80% that she will get a positive mammogram. If a woman does not have breast cancer, the probability is 9.6% that she will also get a positive mammogram. A woman in this age group has a positive mammogram in a routine screening. What is the probability that she actually has breast cancer?"* (Gigerenzer & Hoffrage 1995)
- From this information, we can deduce:
 - $p(C+) = 0.01 \rightarrow p(C-) = 0.99$
 - $p(M+|C+) = 0.8$
 - $p(M+|C-) = 0.096$
- We can now apply Bayes' rule to compute the posterior probability:

$$p(C+ | M+) = \frac{0.8 \cdot 0.01}{0.8 \cdot 0.01 + 0.096 \cdot 0.99} = 0.078$$

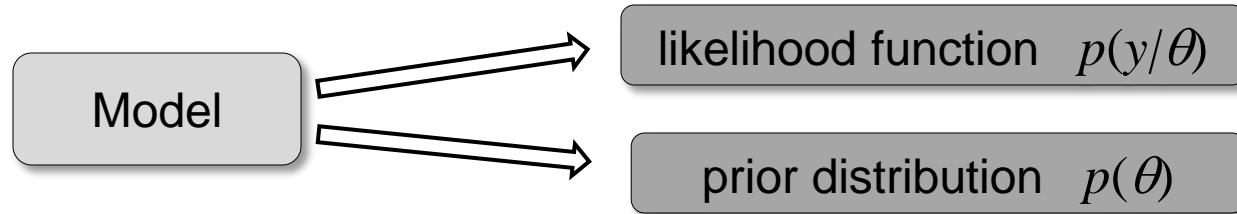
Generative models



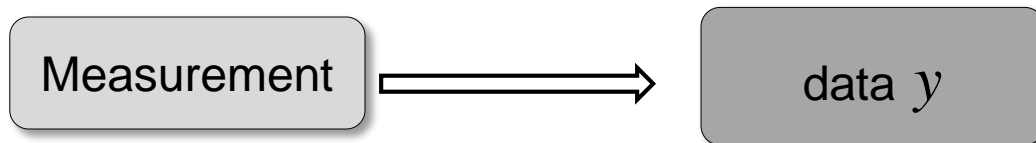
1. specify the joint probability over data (observations) and parameters
2. enforce mechanistic thinking: how could the data have been caused?
3. generate synthetic data (observations) by sampling from the prior – can model explain certain phenomena at all?
4. inference about parameters $\rightarrow p(\theta|y)$
5. model evidence $p(y|m)$: index of model quality

Bayesian inference in practice

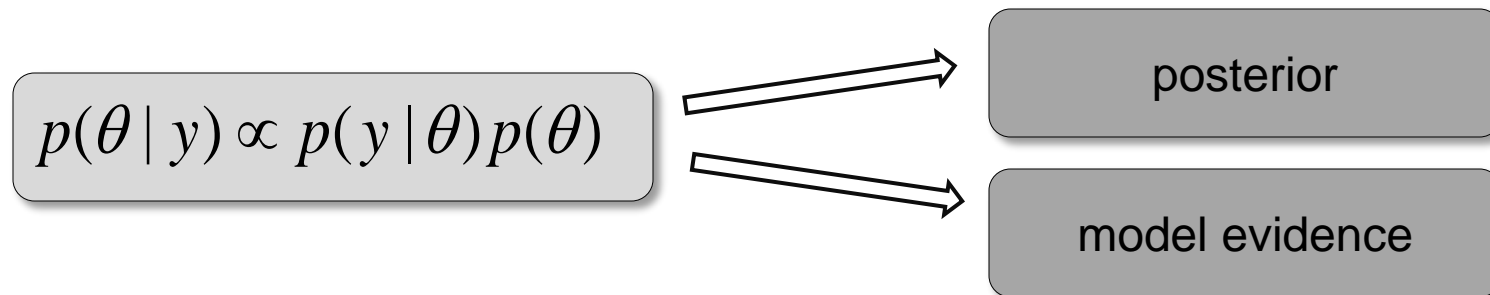
⇒ Formulation of a **generative model**



⇒ Observation of **data**



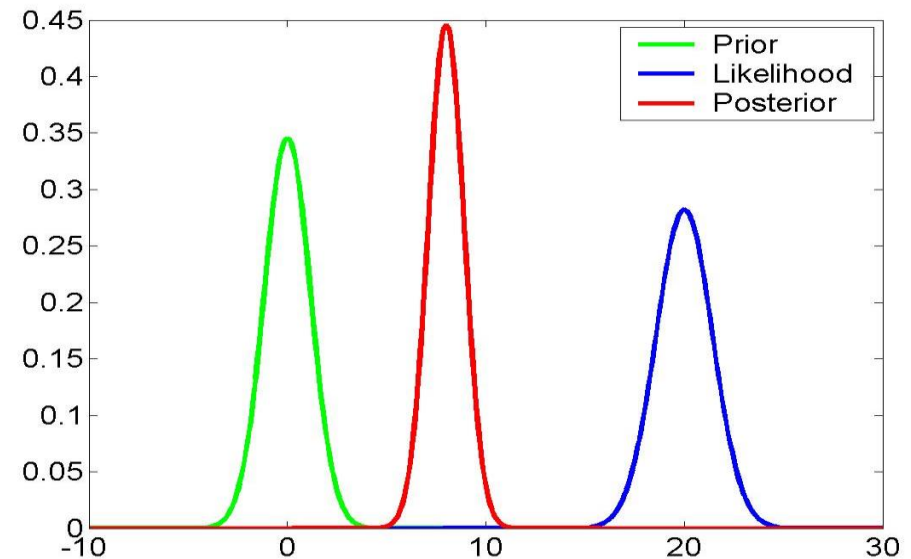
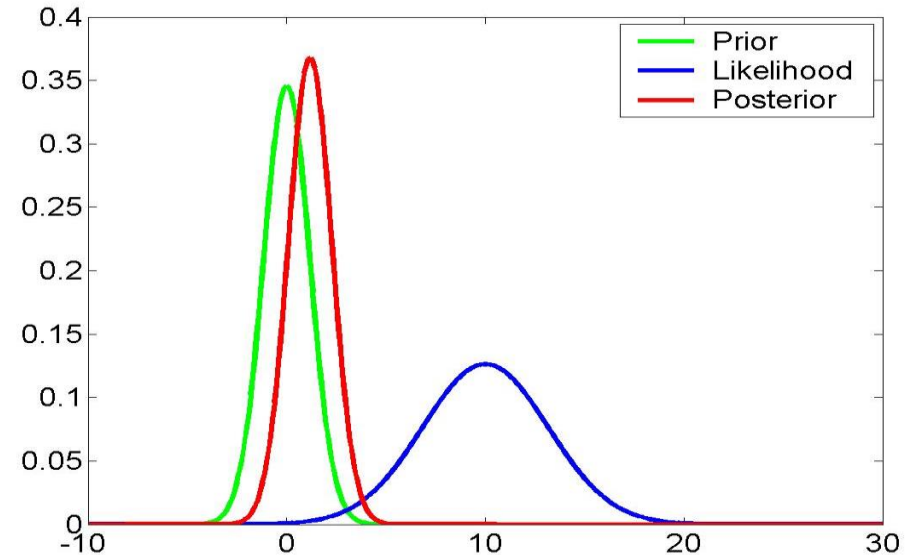
⇒ **Model inversion** – updating one's beliefs



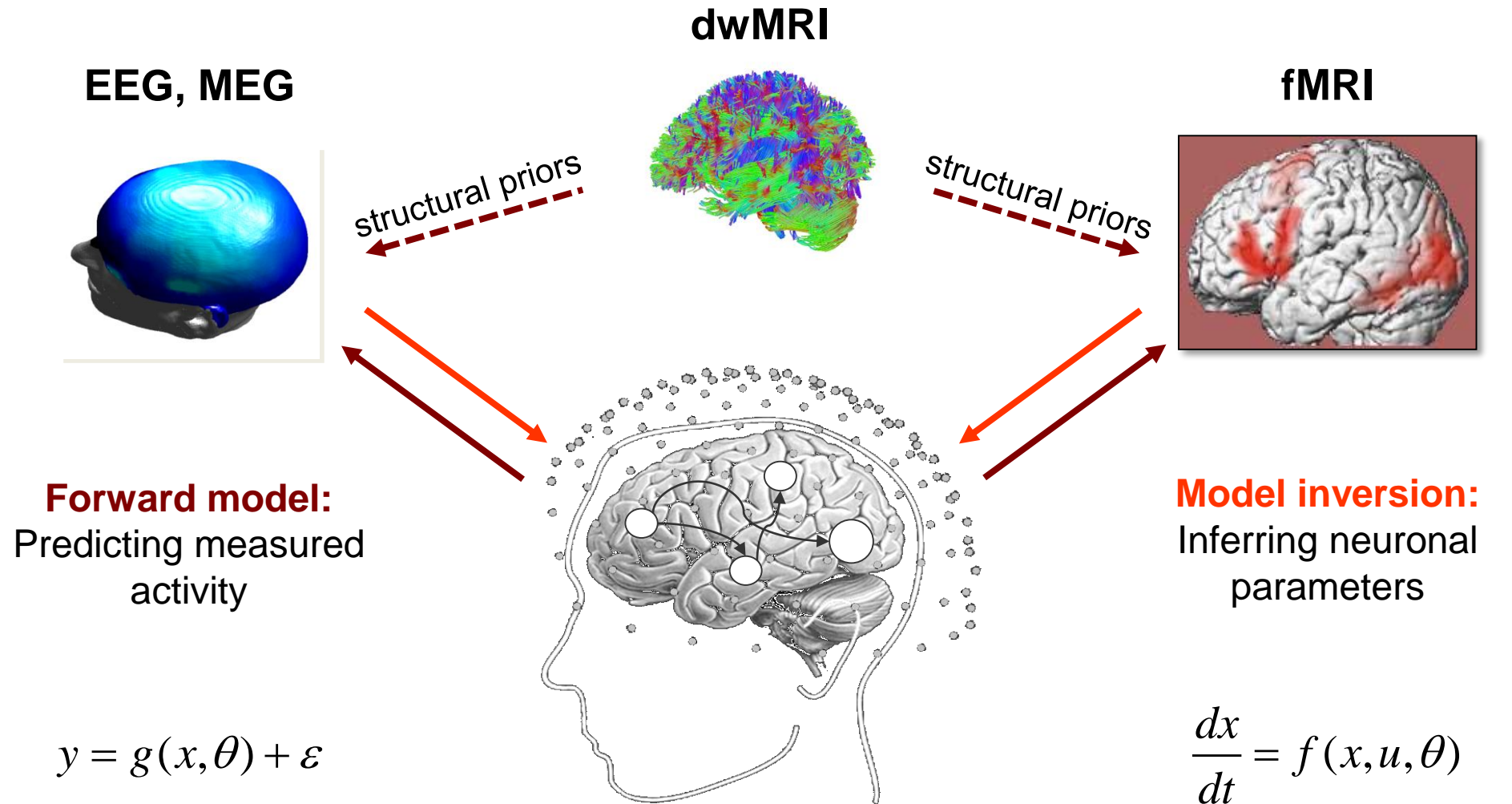
Priors

- Objective priors:
 - "non-informative" priors
 - objective constraints (e.g., non-negativity)
- Subjective priors:
 - subjective but not arbitrary
 - can express beliefs that result from understanding of the problem or system
 - can be result of previous empirical results
- Shrinkage priors:
 - emphasize regularization and sparsity
- Empirical priors:
 - learn parameters of prior distributions from the data ("empirical Bayes")
 - rest on a hierarchical model

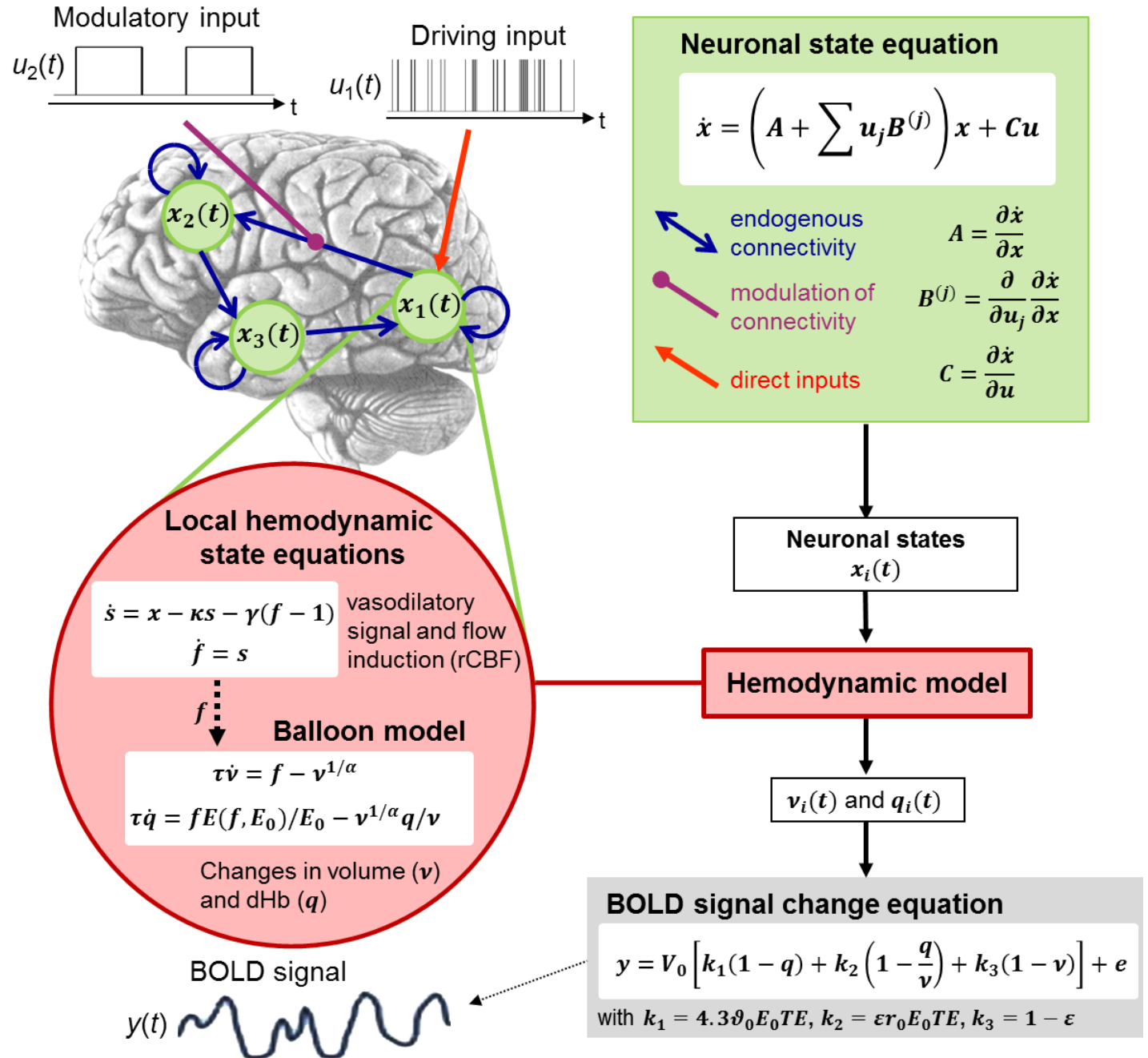
Example of a shrinkage prior

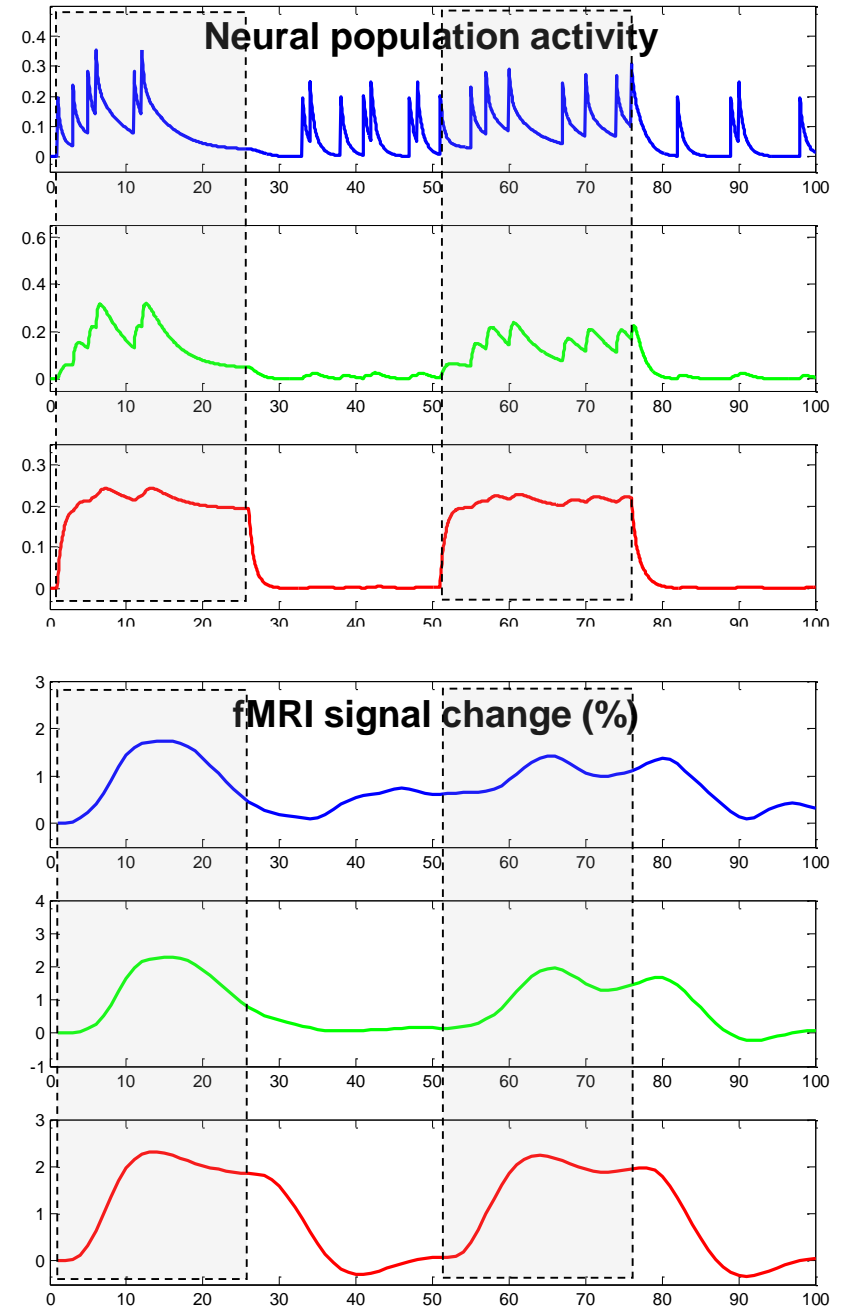
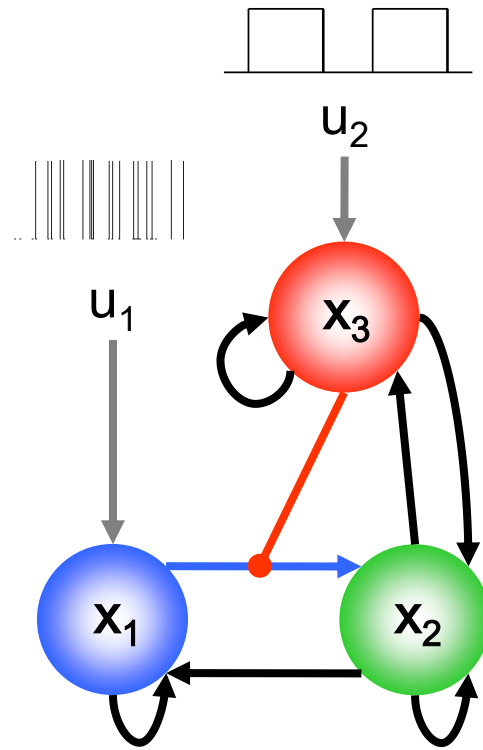
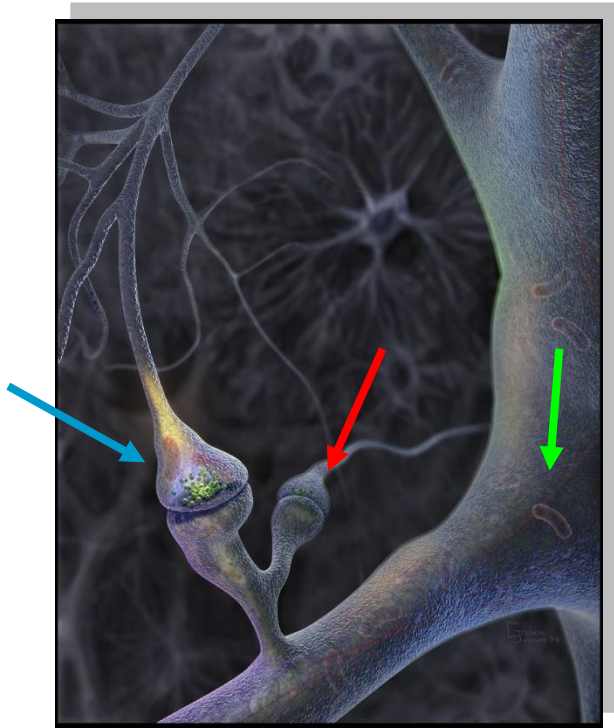


A generative modelling framework for fMRI & EEG: Dynamic causal modeling (DCM)



DCM for fMRI





Nonlinear Dynamic Causal Model for fMRI

$$\frac{dx}{dt} = \left(A + \sum_{i=1}^m u_i B^{(i)} + \sum_{j=1}^n x_j D^{(j)} \right) x + Cu$$

Bayesian system identification

Neural dynamics

$$\frac{dx}{dt} = f(x, u, \theta)$$

Observer function

$$y = g(x, \theta) + \varepsilon$$

$$p(y | \theta, m) = N(g(\theta), \Sigma(\theta))$$

$$p(\theta, m) = N(\mu_\theta, \Sigma_\theta)$$

Inference on model structure

$$p(y | m) = \int p(y | \theta, m) p(\theta) d\theta$$

Inference on parameters

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

$u(t)$



Design experimental inputs

Define likelihood model

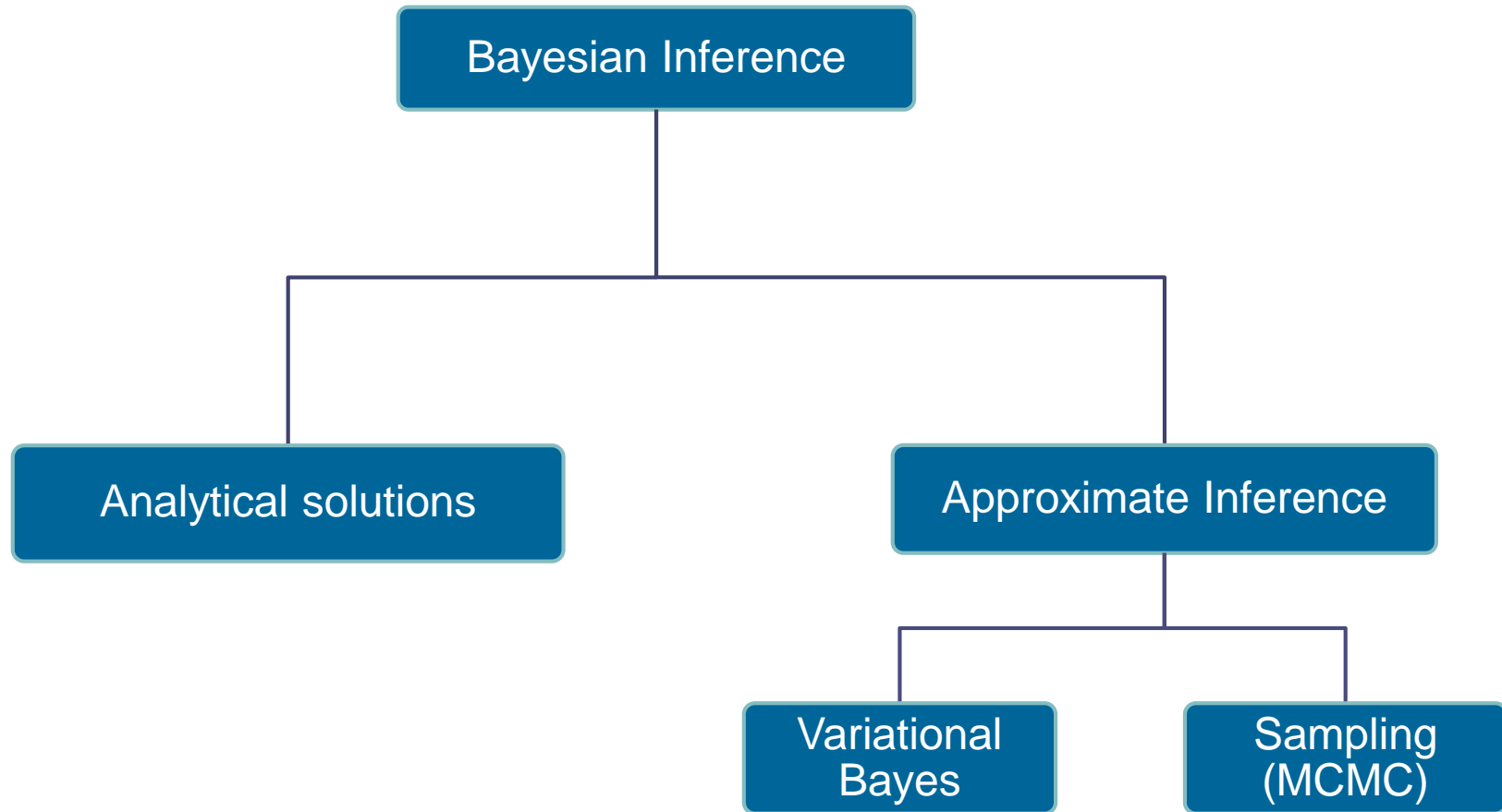
Specify priors

Invert model

Make inferences



Methods for Bayesian inference



How is the posterior computed = how is a generative model inverted?

- **compute the posterior analytically**
 - requires conjugate priors
- **variational Bayes (VB)**
 - often hard work to derive, but fast to compute
 - uses approximations (approximate posteriors, mean field approx.)
 - problems: local minima, potentially inaccurate approximations
- **sampling methods (MCMC)**
 - theoretically guaranteed to be accurate (for infinite computation time)
 - problems: may require very long run time in practice, only heuristics to decide about convergence in practice

Conjugate priors

- for a given likelihood function, the choice of prior determines the algebraic form of the posterior
- for some probability distributions a prior can be found such that the posterior has the same algebraic form as the prior
- such a prior is called “conjugate” to the likelihood
- examples:
 - Normal \propto Normal \times Normal
 - Beta \propto Binomial \times Beta
 - Dirichlet \propto Multinomial \times Dirichlet

$$p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$$

↑ ↑
same form

A simple example: univariate Gaussian belief update

Likelihood & prior

$$p(y | \theta) = N(\theta, \sigma_e^2)$$

$$p(\theta) = N(\mu_{\text{prior}}, \sigma_{\text{prior}}^2)$$

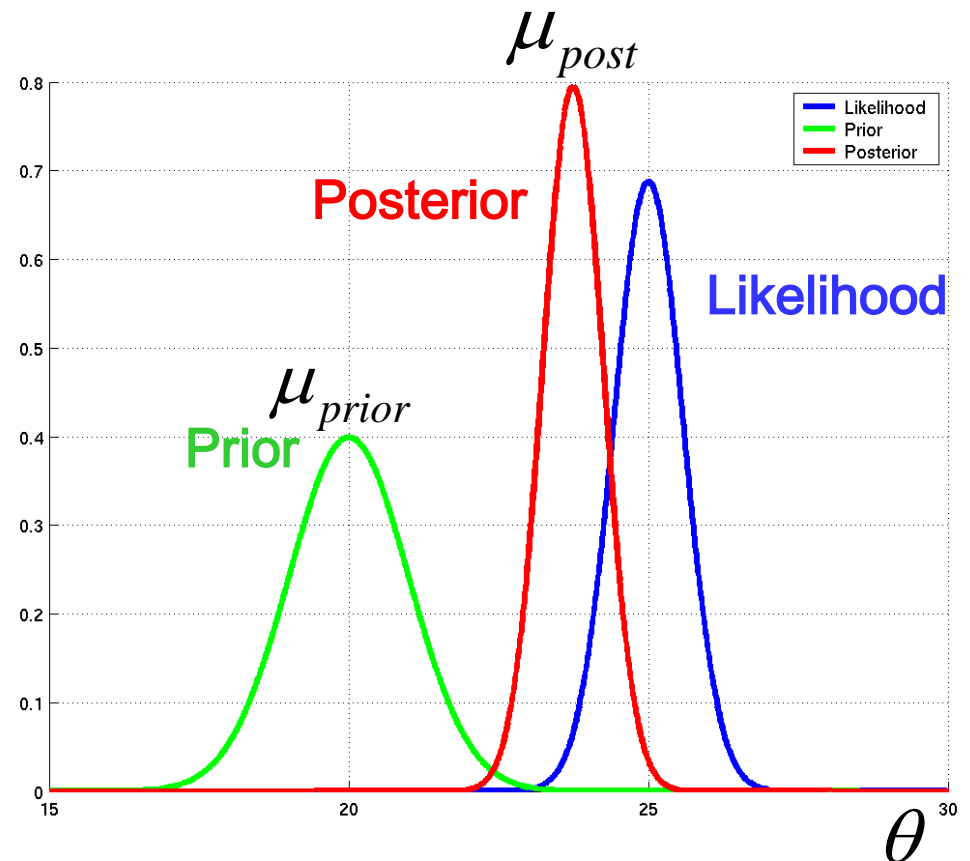
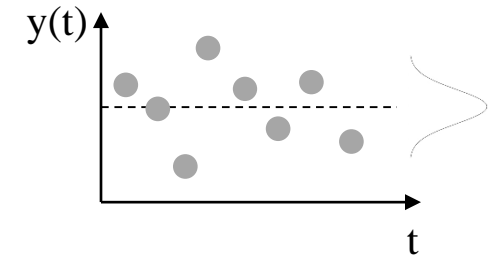
Posterior $p(\theta | y) = N(\mu_{\text{post}}, \lambda_{\text{post}}^{-1})$
(for a single observation y)

$$\frac{1}{\sigma_{\text{post}}^2} = \frac{1}{\sigma_e^2} + \frac{1}{\sigma_{\text{prior}}^2}$$

$$\mu_{\text{post}} = \sigma_{\text{post}}^2 \left(\frac{1}{\sigma_e^2} y + \frac{1}{\sigma_{\text{prior}}^2} \mu_{\text{prior}} \right)$$

posterior mean = variance-weighted combination of prior mean and data

$$y = \theta + \varepsilon$$



Same thing – but expressed as precision weighting

Likelihood & prior

$$p(y | \theta) = N(\theta, \lambda_e^{-1})$$

$$p(\theta) = N(\mu_{\text{prior}}, \lambda_{\text{prior}}^{-1})$$

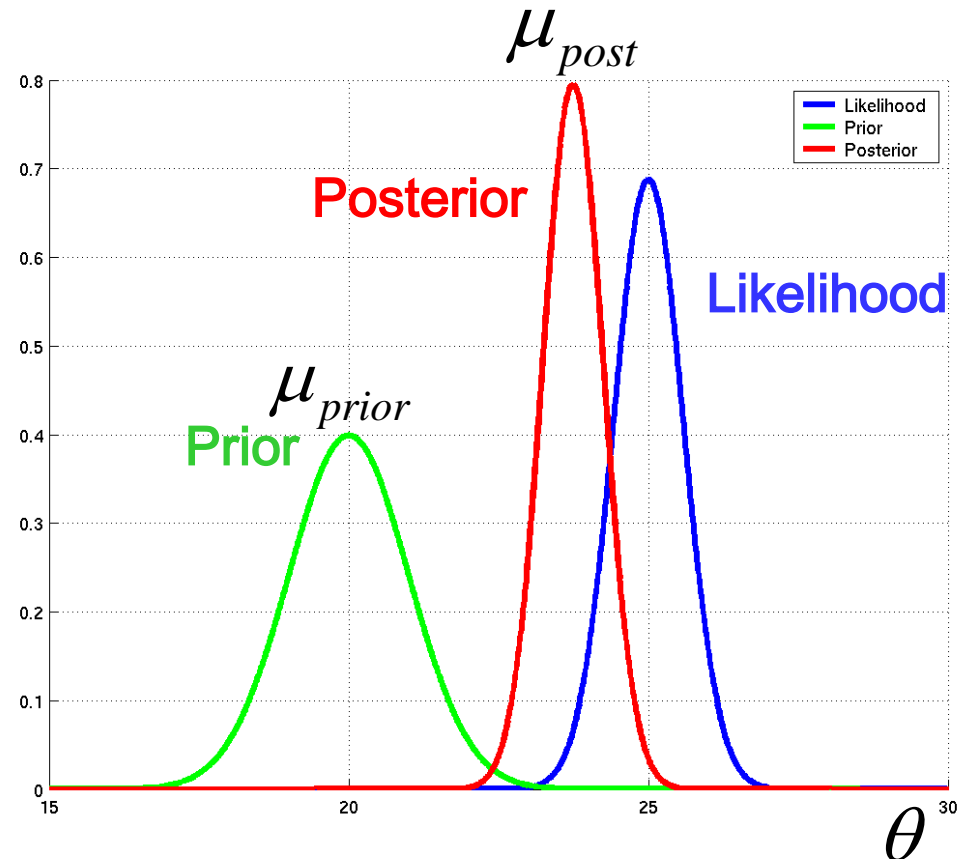
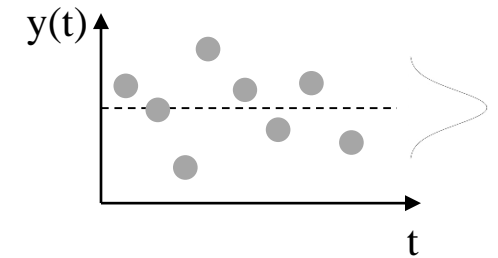
Posterior $p(\theta | y) = N(\mu_{\text{post}}, \lambda_{\text{post}}^{-1})$
(for a single observation y)

$$\lambda_{\text{post}} = \lambda_e + \lambda_{\text{prior}}$$

$$\mu_{\text{post}} = \frac{\lambda_e}{\lambda_{\text{post}}} y + \frac{\lambda_{\text{prior}}}{\lambda_{\text{post}}} \mu_{\text{prior}}$$

relative precision weighting:
posterior mean = precision-weighted
combination of prior mean and data

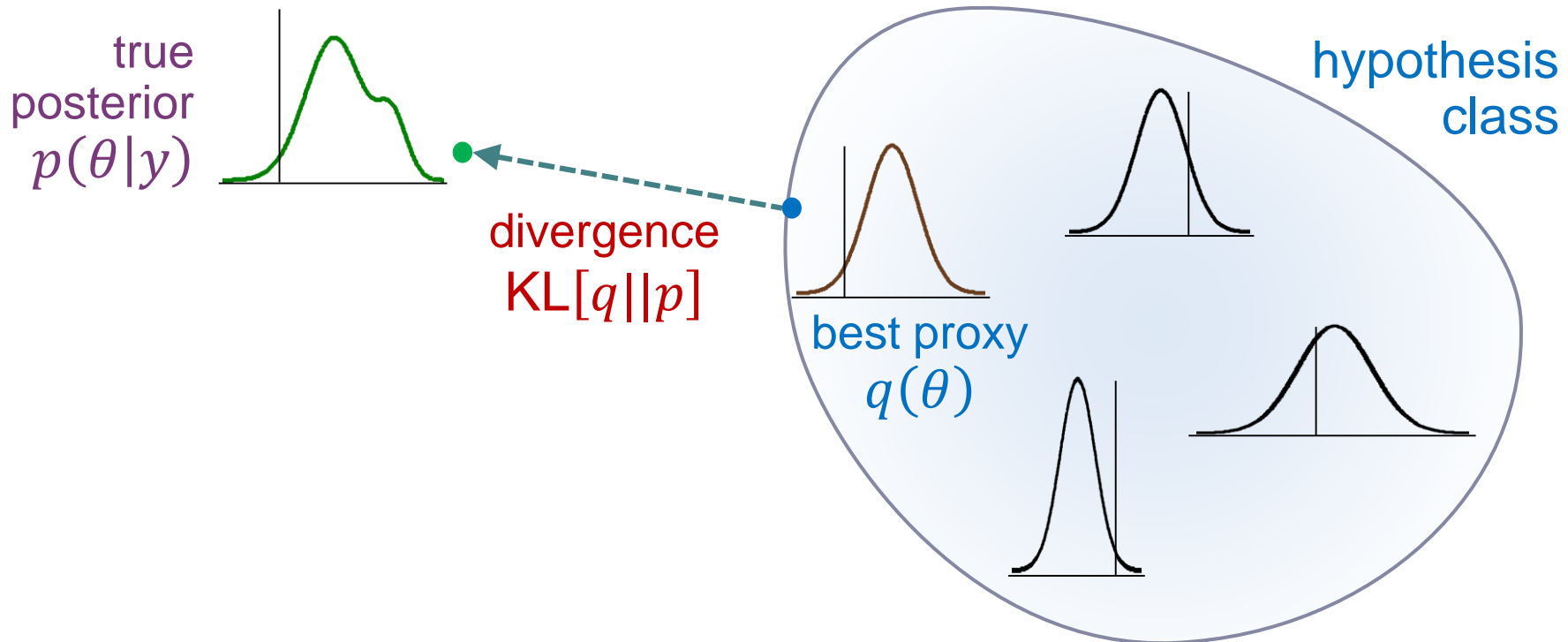
$$y = \theta + \varepsilon$$



Variational Bayes (VB)

Idea: find an approximate density $q(\theta)$ that is maximally similar to the true posterior $p(\theta|y)$.

This is often done by assuming a particular form for q (fixed form VB) and then optimizing its sufficient statistics.



Kullback–Leibler (KL) divergence

- asymmetric measure of the difference between two probability distributions P and Q
- Interpretations of $D_{KL}(P||Q)$:
 - "Bayesian surprise" when Q =prior, P =posterior: measure of the information gained when one updates one's prior beliefs to the posterior P
 - a measure of the information lost when Q is used to approximate P
- non-negative: ≥ 0 (zero when $P=Q$)

$$D_{KL}(P || Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)}$$

$$D_{KL}(P || Q) = \int p(x) \ln \frac{p(x)}{q(x)} dx$$

Variational calculus

Standard calculus

Newton, Leibniz, and others

- functions
 $f: x \mapsto f(x)$
- derivatives $\frac{df}{dx}$

Example: maximize the likelihood expression $p(y|\theta)$ w.r.t. θ

Variational calculus

Euler, Lagrange, and others

- functionals
 $F: f \mapsto F(f)$
- derivatives $\frac{dF}{df}$

Example: maximize the entropy $H[p]$ w.r.t. a probability distribution $p(x)$



Leonhard Euler
(1707 – 1783)

Swiss mathematician,
'Elementa Calculi
Variationum'

Variational Bayes

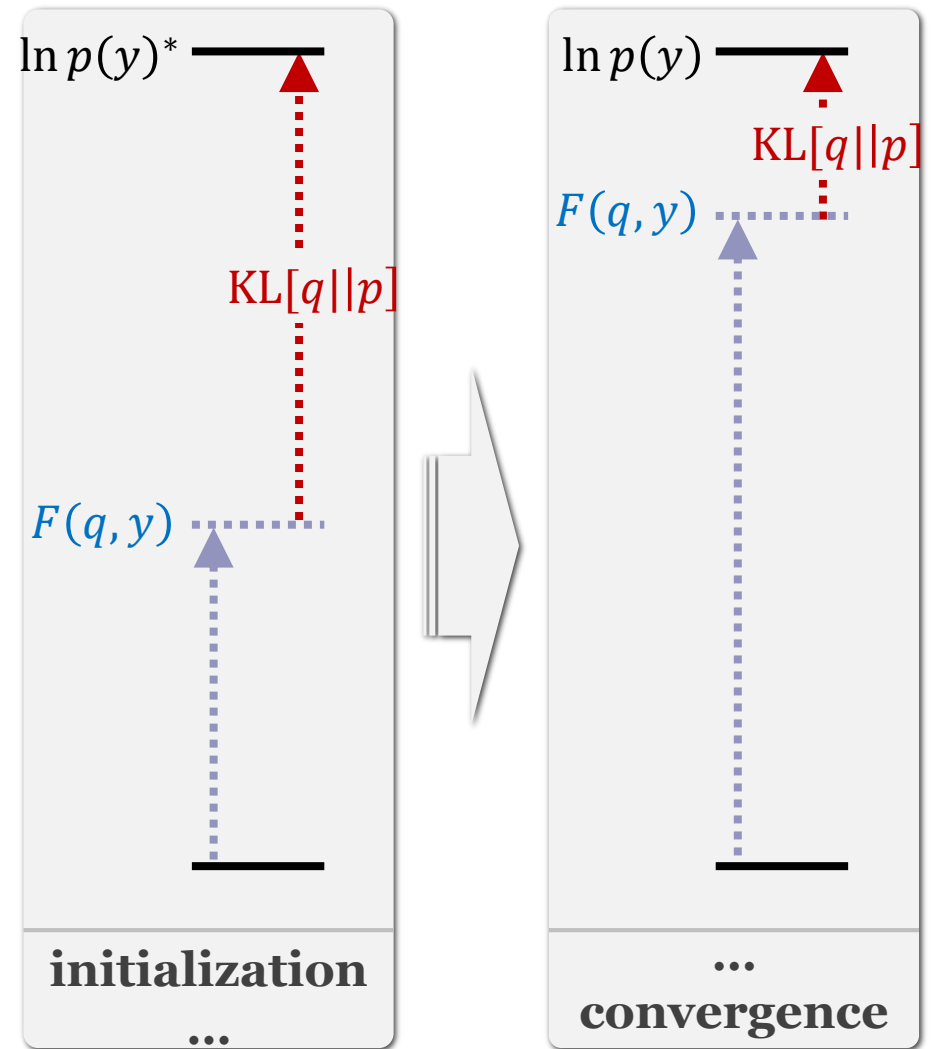
$$\ln p(y) = \underbrace{\text{KL}[q||p]}_{\substack{\text{divergence} \\ \geq 0 \\ \text{(unknown)}}} + \underbrace{F(q, y)}_{\substack{\text{neg. free} \\ \text{energy} \\ \text{(easy to evaluate} \\ \text{for a given } q)}}$$

$F(q)$ is a functional wrt. the approximate posterior $q(\theta)$.

Maximizing $F(q, y)$ is equivalent to:

- minimizing $\text{KL}[q||p]$
- tightening $F(q, y)$ as a lower bound to the log model evidence

When $F(q, y)$ is maximized, $q(\theta)$ is our best estimate of the posterior.



Derivation of the (negative) free energy approximation

- See whiteboard!
- (or Appendix to Stephan et al. 2007, NeuroImage 38: 387-401)

Mean field assumption

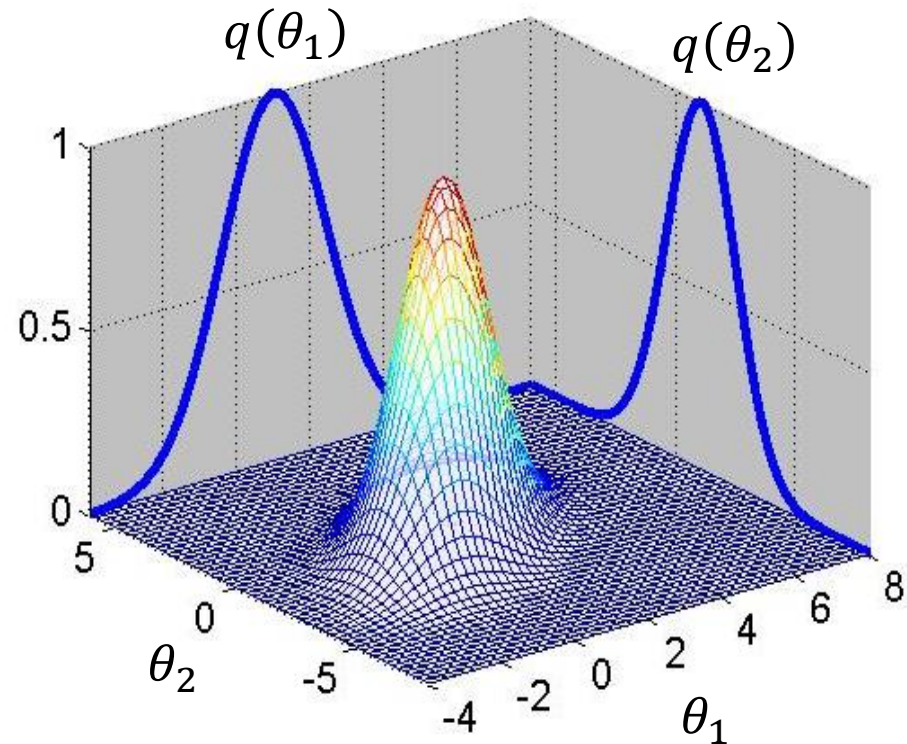
Factorize the approximate posterior $q(\theta)$ into independent partitions:

$$q(\theta) = \prod_i q_i(\theta_i)$$

where $q_i(\theta_i)$ is the approximate posterior for the i^{th} subset of parameters.

For example, split parameters and hyperparameters:

$$p(\theta, \lambda | y) \approx q(\theta, \lambda) = q(\theta)q(\lambda)$$



VB in a nutshell (under mean-field approximation)

- 1 Neg. free-energy approx. to model evidence.

$$\ln p(y|m) = F + KL[q(\theta, \lambda), p(\theta, \lambda | y)]$$

$$F = \langle \ln p(y | \theta, \lambda) \rangle_q - KL[q(\theta, \lambda), p(\theta, \lambda | m)]$$

- 2 Mean field approx.

$$p(\theta, \lambda | y) \approx q(\theta, \lambda) = q(\theta)q(\lambda)$$

- 3 Maximise neg. free energy wrt. q = minimise divergence, by maximising variational energies

$$q(\theta) \propto \exp(I_\theta) = \exp\left[\langle \ln p(y, \theta, \lambda) \rangle_{q(\lambda)}\right]$$

$$q(\lambda) \propto \exp(I_\lambda) = \exp\left[\langle \ln p(y, \theta, \lambda) \rangle_{q(\theta)}\right]$$

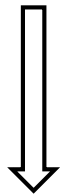
- 4 Iterative updating of sufficient statistics of approx. posteriors by gradient ascent.

Model comparison and selection

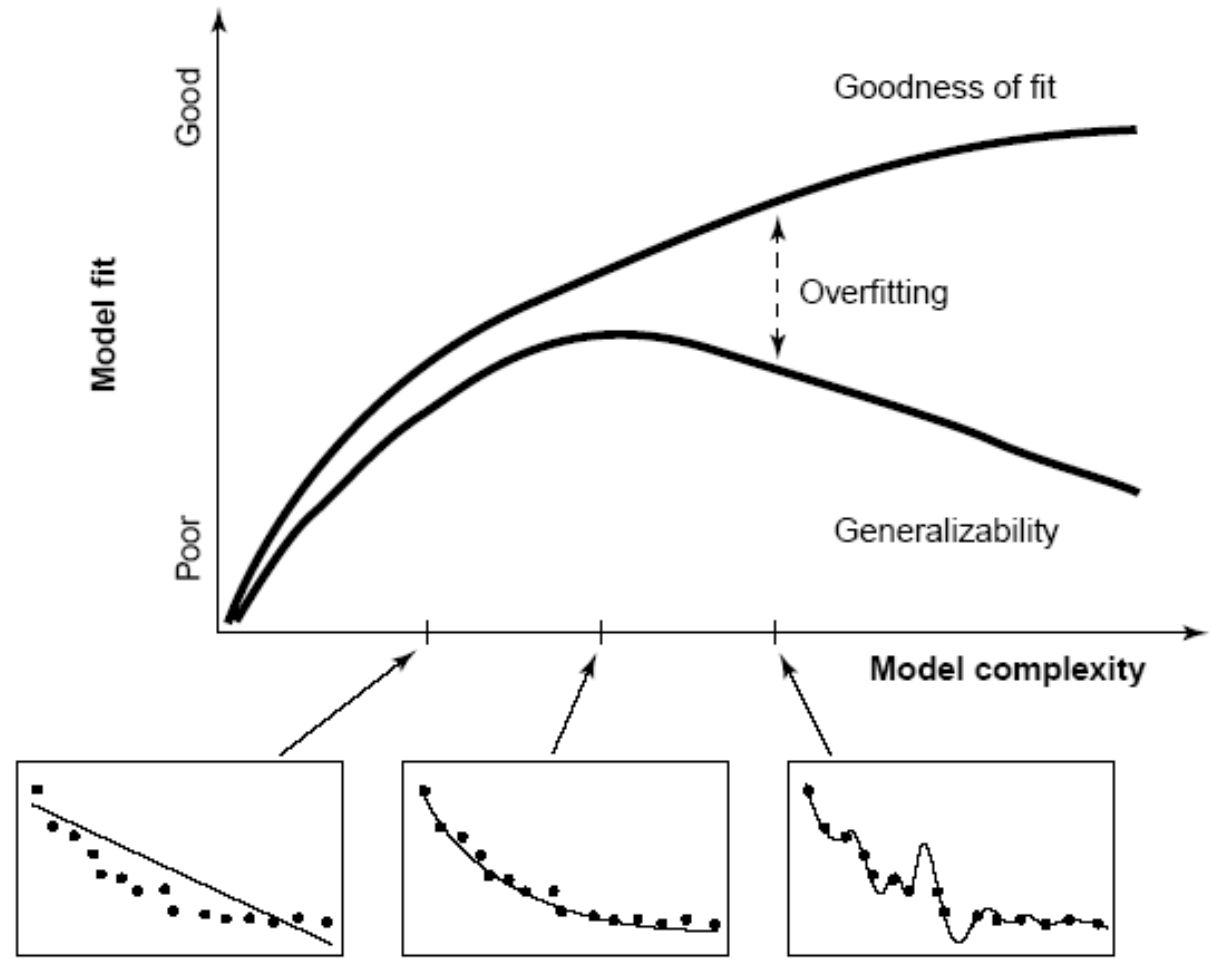
Given competing hypotheses on structure & functional mechanisms of a system, which model is the best?



Which model represents the best balance between model fit and model complexity?



For which model m does $p(y|m)$ become maximal?



Bayesian model selection (BMS)

- First step of inference: define model space M

$$|M| \in [1, \infty[$$

- Inference on model structure m :

Posterior model probability

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

- For a uniform prior on m , model evidence sufficient for model selection

Model evidence:

$$p(y | m) = \int p(y | \theta, m) p(\theta | m) d\theta$$

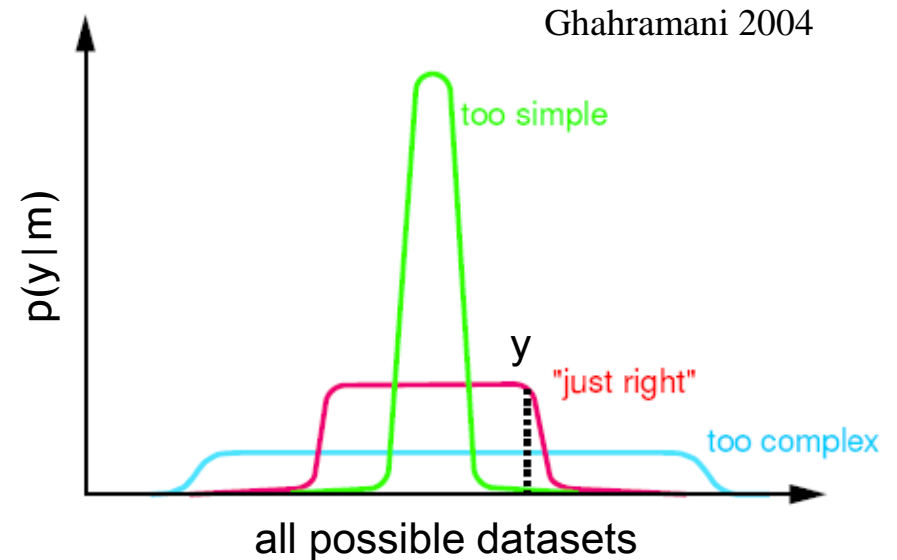
Bayesian model selection (BMS)

Model evidence:

$$p(y | m) = \int p(y | \theta, m) p(\theta | m) d\theta$$

⇒ probability that data were generated by model m , averaging over all possible parameter values (with probability weights as specified by the prior)

⇒ accounts for both accuracy and complexity of the model



Various approximations:

- negative free energy (F)
- Akaike Information Criterion (AIC)
- Bayesian Information Criterion (BIC)

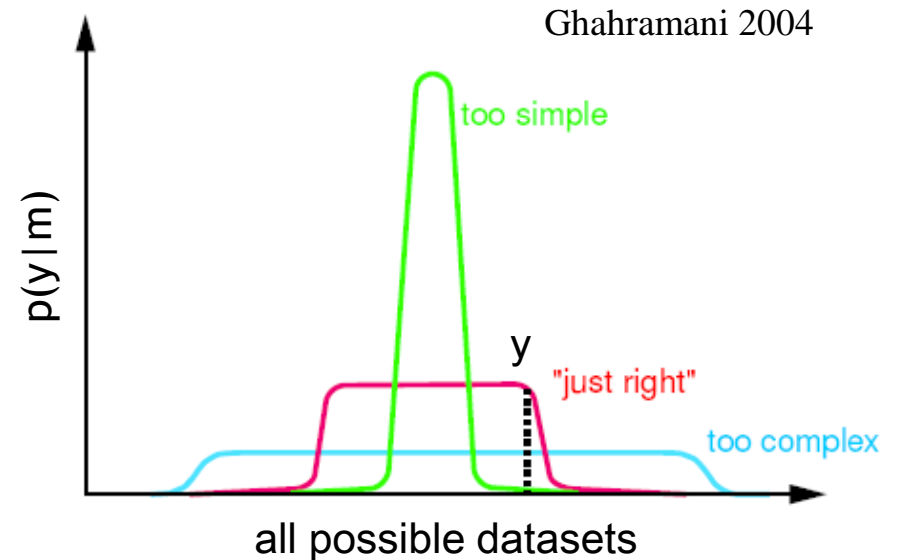
Bayesian model selection (BMS)

Model evidence:

$$p(y | m) = \int p(y | \theta, m) p(\theta | m) d\theta$$

⇒ “If I randomly sampled from my prior and plugged the resulting value into the likelihood function, how close would the predicted data be – on average – to my observed data?”

⇒ accounts for both accuracy and complexity of the model



Various approximations:

- negative free energy (F)
- Akaike Information Criterion (AIC)
- Bayesian Information Criterion (BIC)

Approximations to the model evidence

Logarithm is a
monotonic function



Maximizing log model evidence
= Maximizing model evidence

Log model evidence = balance between fit and complexity

$$\begin{aligned}\log p(y | m) &= \textit{accuracy}(m) - \textit{complexity}(m) \\ &= \log p(y | \theta, m) - \textit{complexity}(m)\end{aligned}$$

Akaike Information Criterion:

$$AIC = \log p(y | \theta, m) - p$$

No. of
parameters

No. of
data points

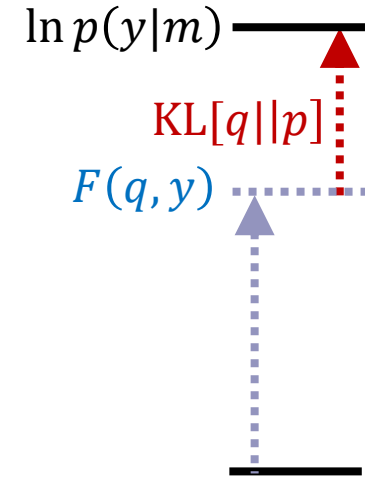
Bayesian Information Criterion:

$$BIC = \log p(y | \theta, m) - \frac{p}{2} \log N$$

The (negative) free energy approximation F

F is a lower bound on the log model evidence:

$$\log p(y | m) = F + KL[q(\theta), p(\theta | y, m)]$$



Like AIC/BIC, F is an accuracy/complexity tradeoff:

$$F = \underbrace{\langle \log p(y | \theta, m) \rangle}_{\text{accuracy}} - \underbrace{KL[q(\theta), p(\theta | m)]}_{\text{complexity}}$$

The (negative) free energy approximation

- Log evidence is thus expected log likelihood (wrt. q) plus 2 KL's:

$$\log p(y | m)$$

$$= \langle \log p(y | \theta, m) \rangle - KL[q(\theta), p(\theta | m)] + KL[q(\theta), p(\theta | y, m)]$$

$$F = \log p(y | m) - KL[q(\theta), p(\theta | y, m)]$$

$$= \underbrace{\langle \log p(y | \theta, m) \rangle}_{\text{accuracy}} - \underbrace{KL[q(\theta), p(\theta | m)]}_{\text{complexity}}$$

accuracy

complexity

The complexity term in F

- In contrast to AIC & BIC, the complexity term of the negative free energy F accounts for parameter interdependencies.

Under Gaussian assumptions about the posterior (Laplace approximation):

$$\begin{aligned} & KL[q(\theta), p(\theta | m)] \\ &= \frac{1}{2} \ln |C_\theta| - \frac{1}{2} \ln |C_{\theta|y}| + \frac{1}{2} (\mu_{\theta|y} - \mu_\theta)^T C_\theta^{-1} (\mu_{\theta|y} - \mu_\theta) \end{aligned}$$

- The complexity term of F is higher
 - the more independent the prior parameters (\uparrow effective DFs)
 - the more dependent the posterior parameters
 - the more the posterior mean deviates from the prior mean

Bayes factors

To compare two models, we could just compare their log evidences.

But: the log evidence is just some number – not very intuitive!

A more intuitive interpretation of model comparisons is made possible by Bayes factors:

$$B_{12} = \frac{p(y | m_1)}{p(y | m_2)}$$

positive value, $[0; \infty[$

Kass & Raftery classification:

B_{12}	$p(m_1 y)$	Evidence
1 to 3	50-75%	weak
3 to 20	75-95%	positive
20 to 150	95-99%	strong
≥ 150	$\geq 99\%$	Very strong

Fixed effects BMS at group level

Group Bayes factor (GBF) for $1 \dots K$ subjects:

$$GBF_{ij} = \prod_k BF_{ij}^{(k)}$$

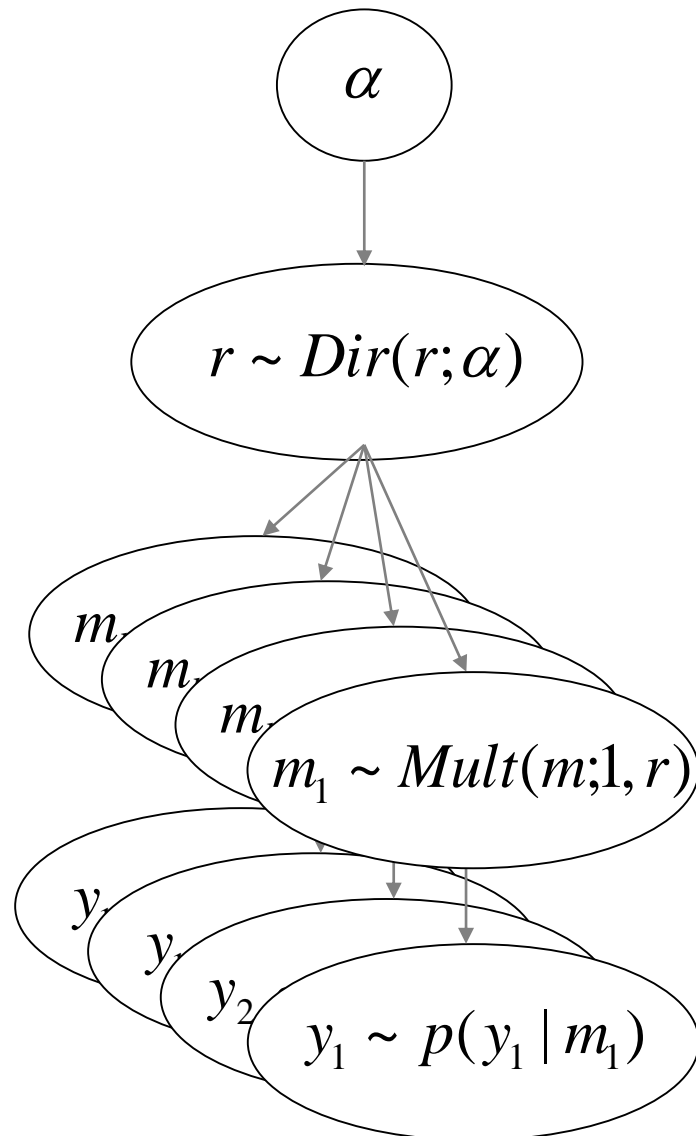
Average Bayes factor (ABF):

$$ABF_{ij} = \sqrt[K]{\prod_k BF_{ij}^{(k)}}$$

Problems:

- blind with regard to group heterogeneity
- sensitive to outliers

Random effects BMS for heterogeneous groups



Dirichlet parameters α
= “occurrences” of models in the population

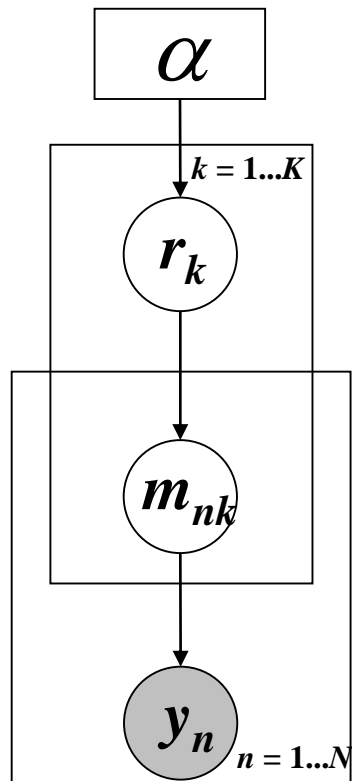
Dirichlet distribution of model probabilities r

Multinomial distribution of model labels m

Measured data y

**Model inversion
by Variational
Bayes or MCMC**

Random effects BMS for heterogeneous groups



Dirichlet parameters α
= “occurrences” of models in the population

Dirichlet distribution of model probabilities r

Multinomial distribution of model labels m

Measured data y

**Model inversion
by Variational
Bayes or MCMC**

Four equivalent options for reporting model ranking by random effects BMS

1. **Dirichlet parameter estimates**

α

2. **expected posterior probability** of obtaining the k -th model for any randomly selected subject

$$\langle r_k \rangle_q = \alpha_k / (\alpha_1 + \dots + \alpha_K)$$

3. **exceedance probability** that a particular model k is more likely than any other model (of the K models tested), given the group data

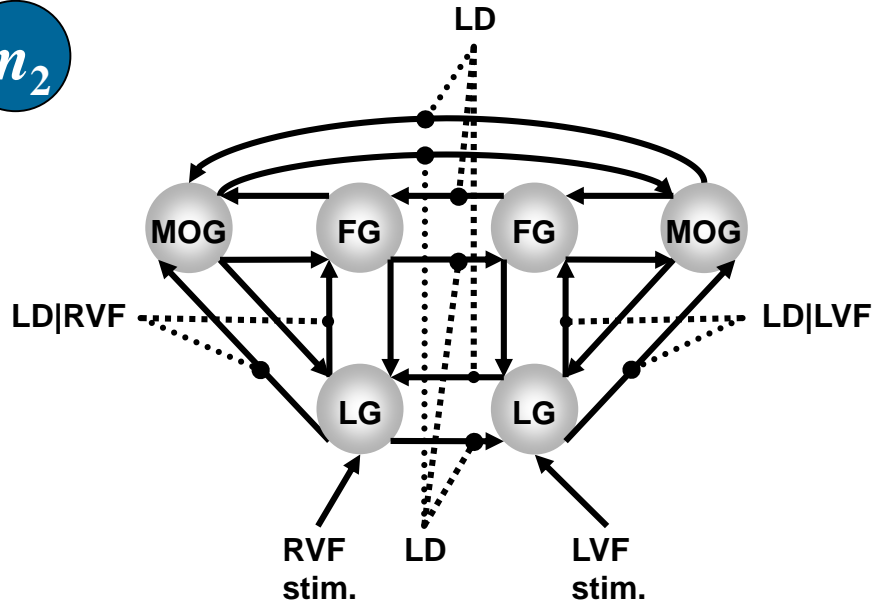
$$\exists k \in \{1 \dots K\}, \forall j \in \{1 \dots K \mid j \neq k\} :$$

$$\varphi_k = p(r_k > r_j \mid y; \alpha)$$

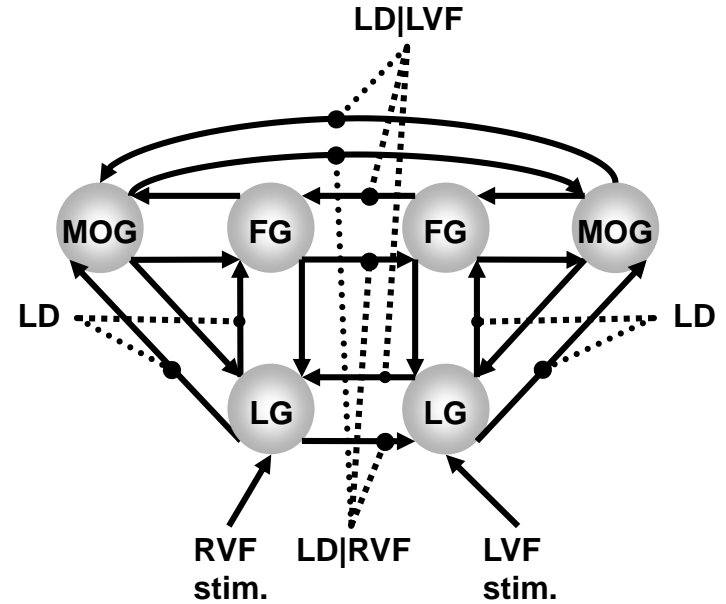
4. **protected exceedance probability:**
see below

Example: Hemispheric interactions during vision

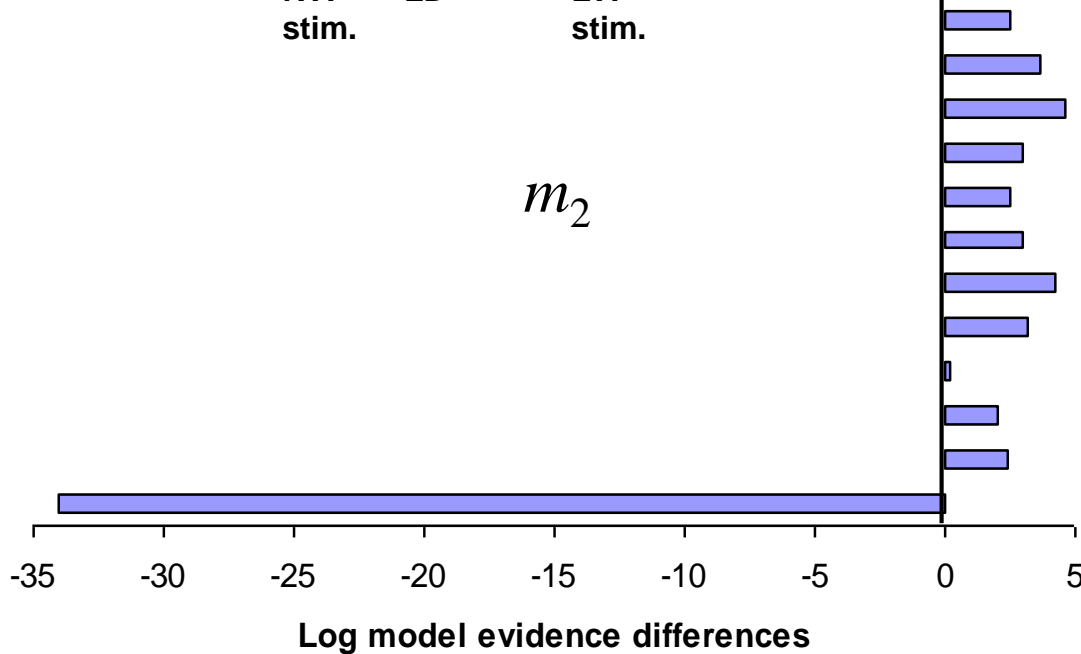
m_2



m_1

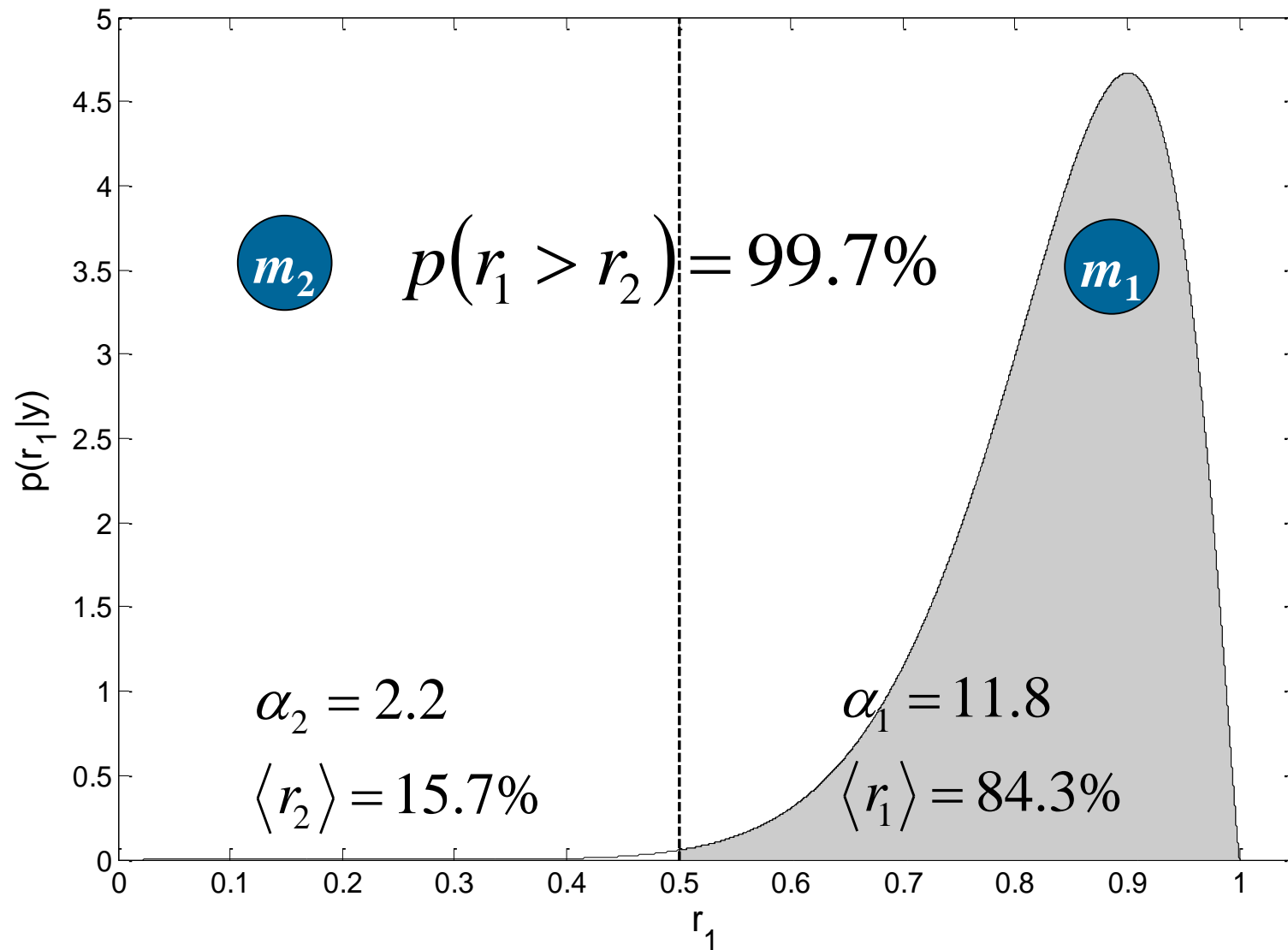


Subjects



m_1

Data: Stephan et al. 2003, *Science*
Models: Stephan et al. 2007, *J. Neurosci.*



Example: Synaesthesia

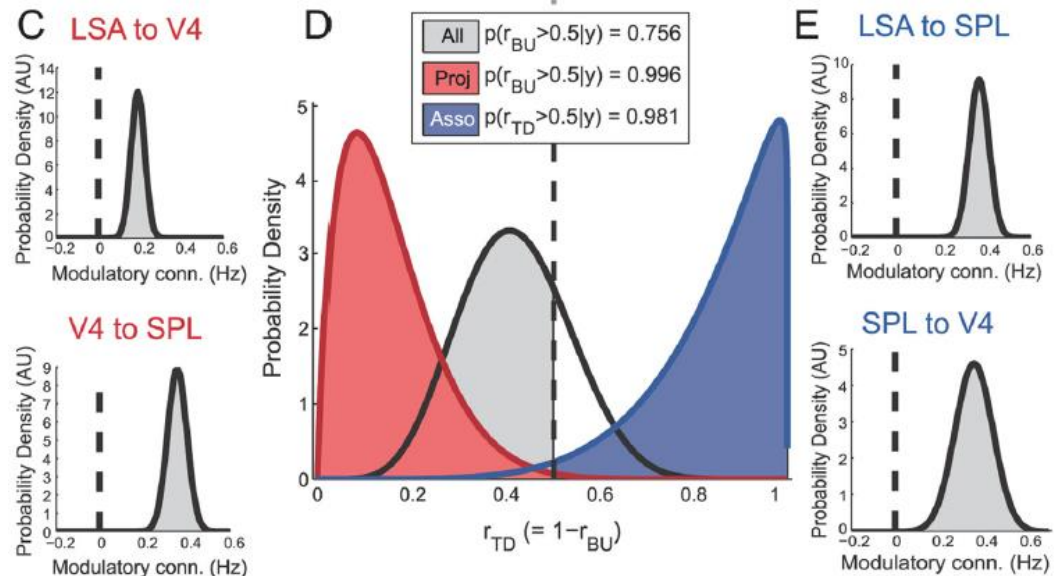
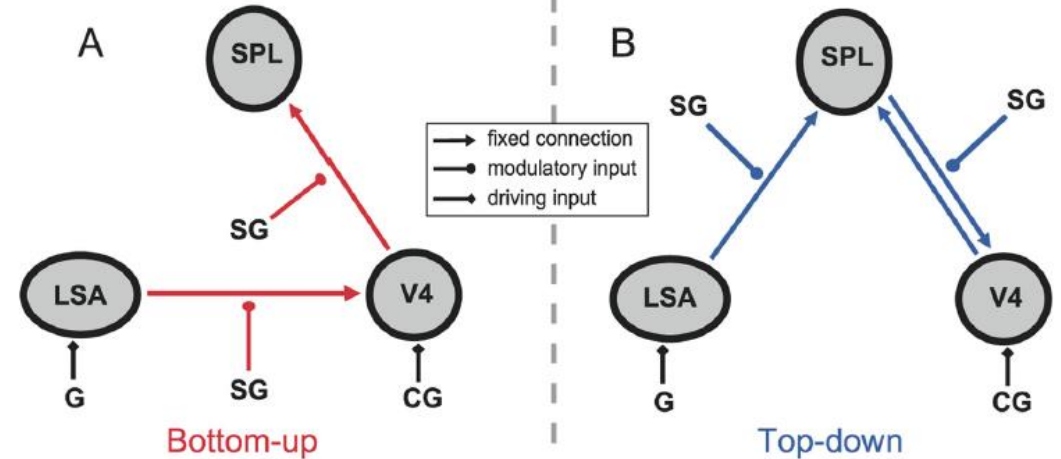
- “projectors” experience color externally colocalized with a presented grapheme
- “associators” report an internally evoked association
- across all subjects: no evidence for either model
- but BMS results map precisely onto projectors (bottom-up mechanisms) and associators (top-down)

PROJECTORS

AB

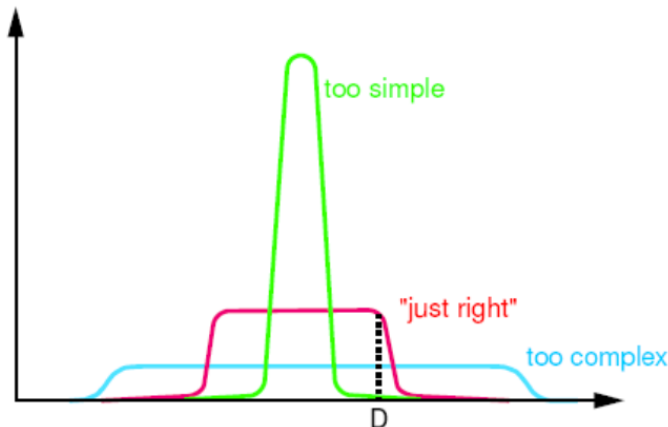
ASSOCIATORS

AB



Overfitting at the level of models

- \uparrow #models \Rightarrow \uparrow risk of overfitting
- solutions:
 - regularisation: definition of model space = choosing priors $p(m)$
 - family-level BMS
 - Bayesian model averaging (BMA)



posterior model probability:

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

BMA:

$$p(\theta | y) = \sum_m p(\theta | y, m) p(m | y)$$

Model space partitioning: comparing model families

- partitioning model space into K subsets or families:

$$M = \{f_1, \dots, f_K\}$$

- pooling information over all models in these subsets allows one to compute the probability of a model family, given the data

$$p(f_k)$$

- effectively removes uncertainty about any aspect of model structure, other than the attribute of interest (which defines the partition)

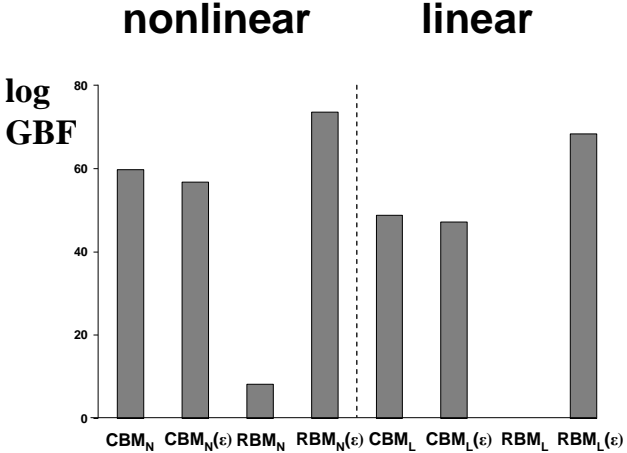
Family-level inference: random effects – a special case

- When the families are of equal size, one can simply sum the posterior model probabilities within families by exploiting the agglomerative property of the Dirichlet distribution:

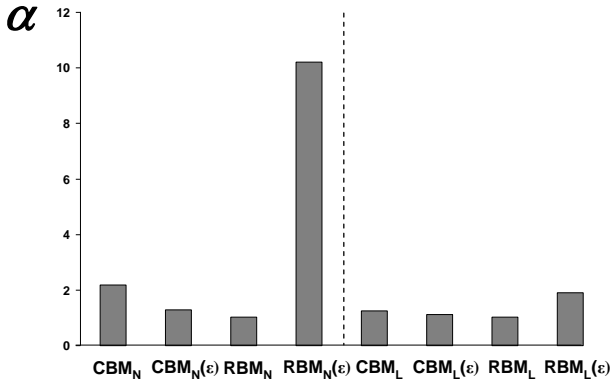
$$\begin{aligned} (r_1, r_2, \dots, r_K) &\sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_K) \\ \Rightarrow r_1^* &= \sum_{k \in N_1} r_k, r_2^* = \sum_{k \in N_2} r_k, \dots, r_J^* = \sum_{k \in N_J} r_k \\ &\sim \text{Dir} \left(\alpha_1^* = \sum_{k \in N_1} \alpha_k, \alpha_2^* = \sum_{k \in N_2} \alpha_k, \dots, \alpha_J^* = \sum_{k \in N_J} \alpha_k \right) \end{aligned}$$

Model space partitioning: comparing model families

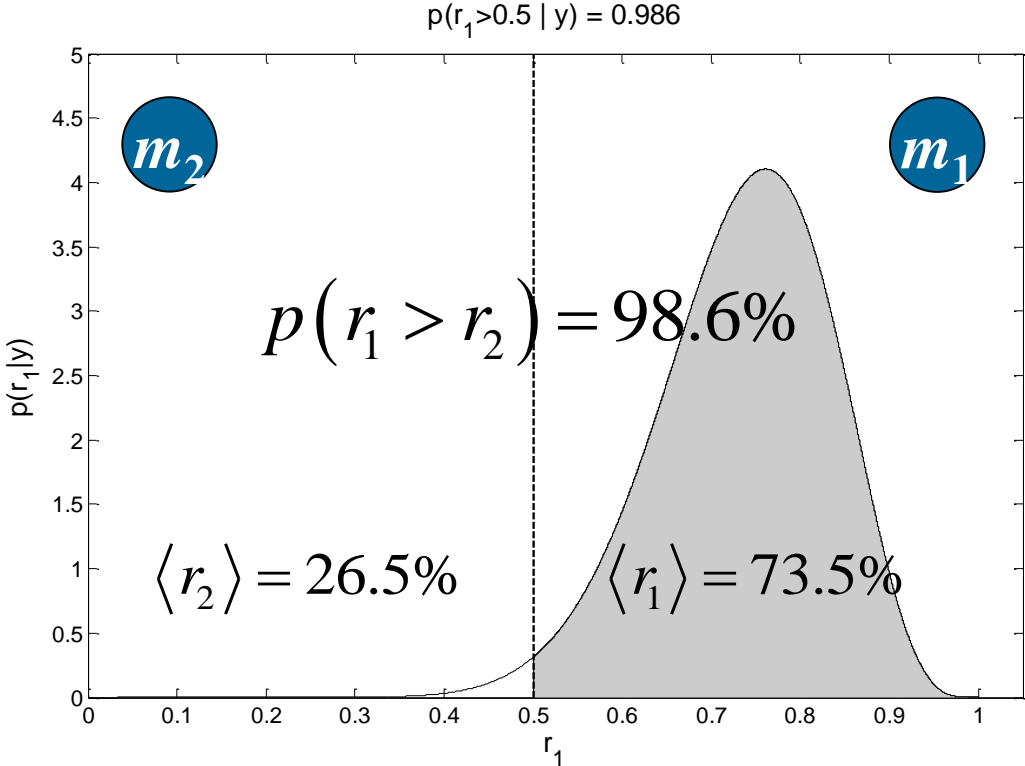
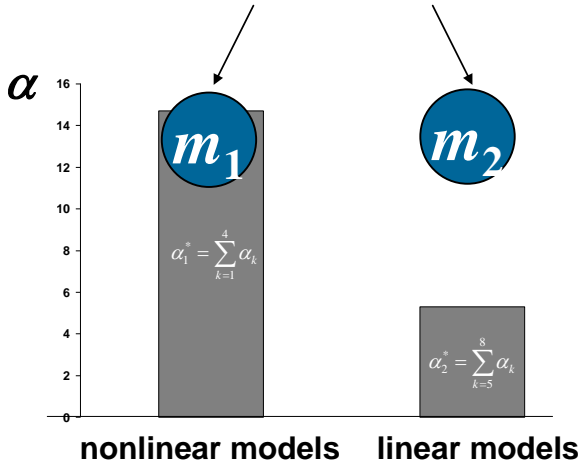
FFX



RFX

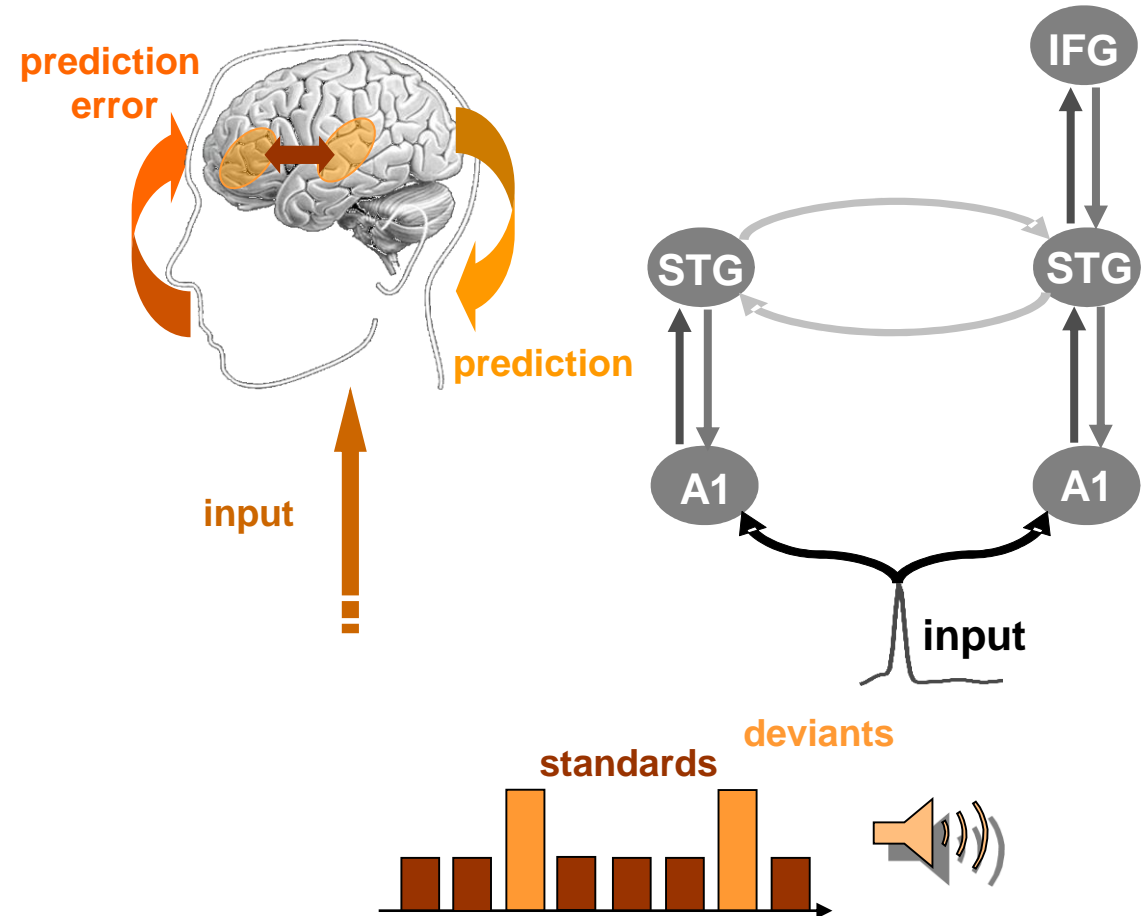
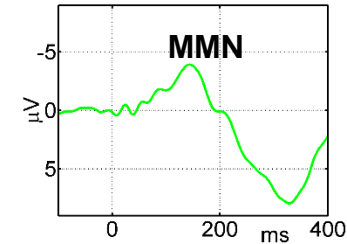


Model space partitioning

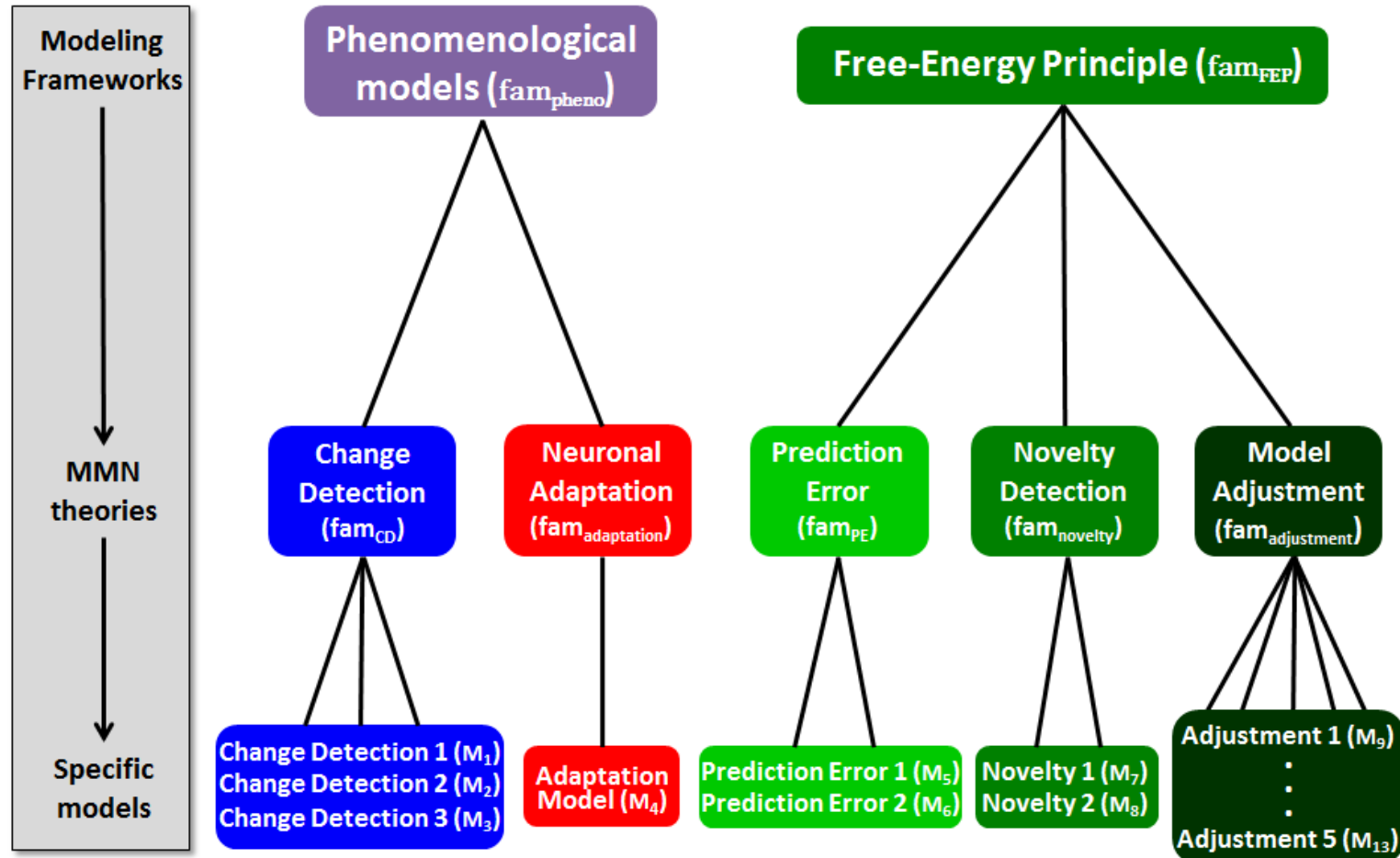


Mismatch negativity (MMN)

- elicited by surprising stimuli (scales with unpredictability)
- ↓ in schizophrenic patients
- classical interpretations:
 - pre-attentive change detection
 - neuronal adaptation
- current theories:
 - reflection of (hierarchical) Bayesian inference



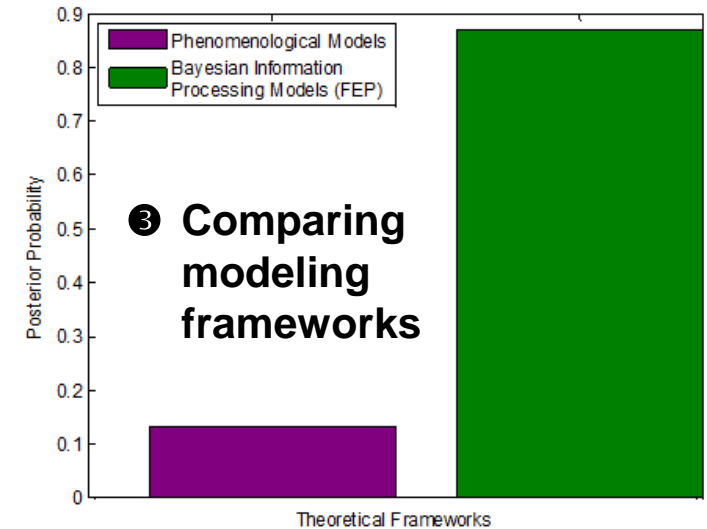
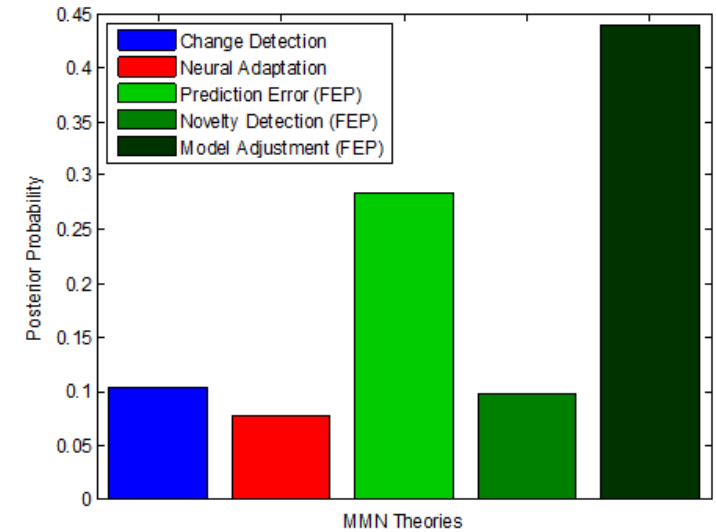
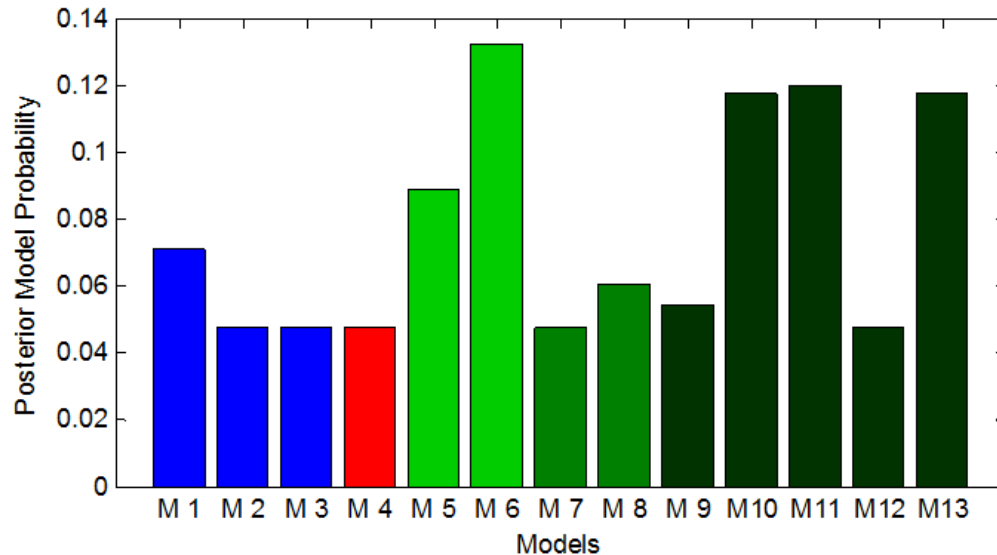
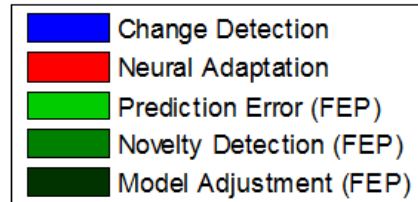
Modelling Trial-by-Trial Changes of the Mismatch Negativity (MMN)



MMN model comparison at multiple levels

② Comparing MMN theories

① Comparing individual models



Bayesian Model Averaging (BMA)

- abandons dependence of parameter inference on a single model and takes into account model uncertainty
- uses the entire model space considered (or an optimal family of models)
- averages parameter estimates, weighted by posterior model probabilities
- represents a particularly useful alternative
 - when none of the models (or model subspaces) considered clearly outperforms all others
 - when comparing groups for which the optimal model differs

single-subject BMA:

$$p(\theta | y) \\ = \sum_m p(\theta | y, m) p(m | y)$$

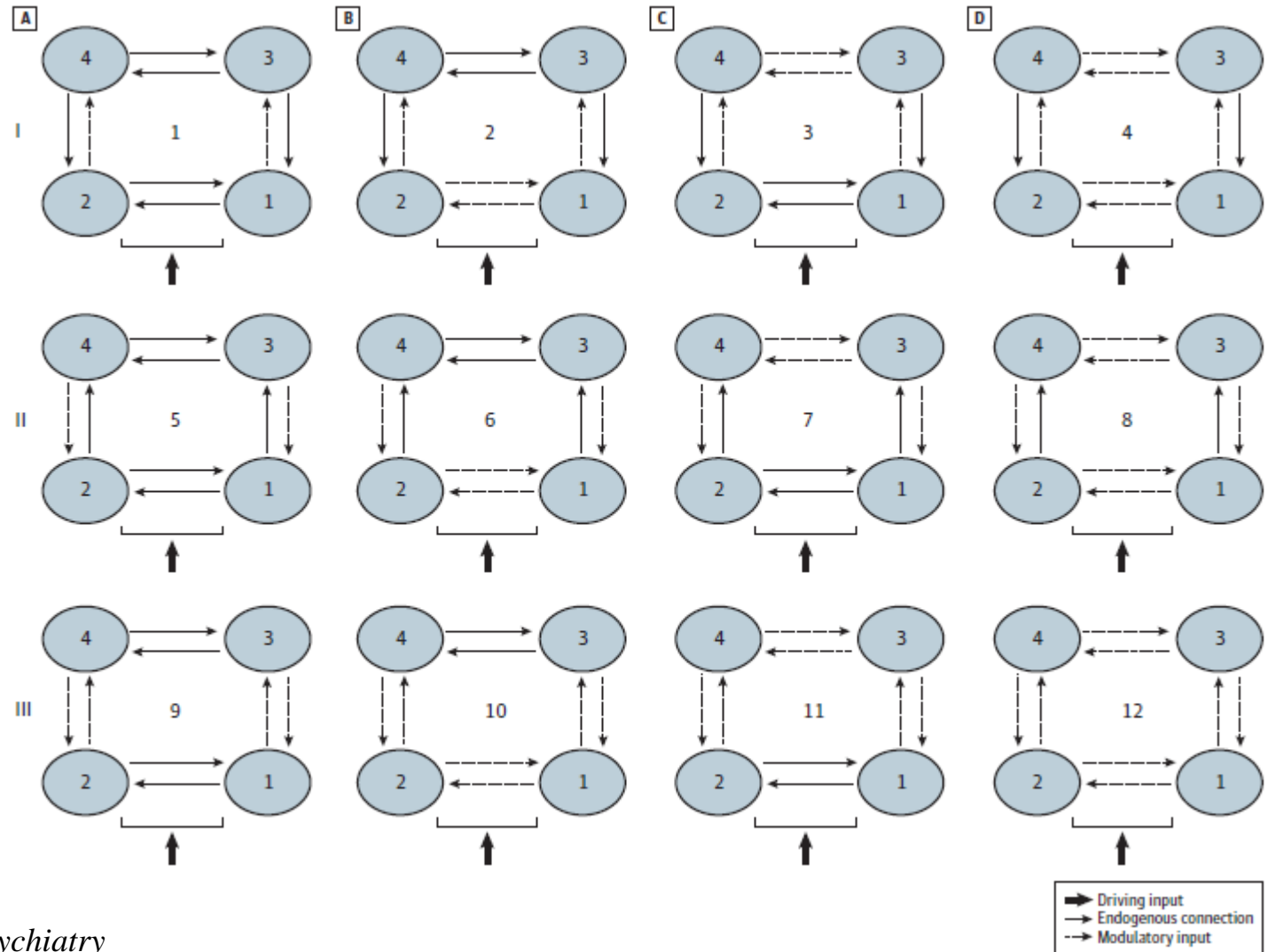
group-level BMA:

$$p(\theta_n | y_{1..N}) \\ = \sum_m p(\theta_n | y_n, m) p(m | y_{1..N})$$

NB: $p(m|y_{1..N})$ can be obtained by either FFX or RFX BMS

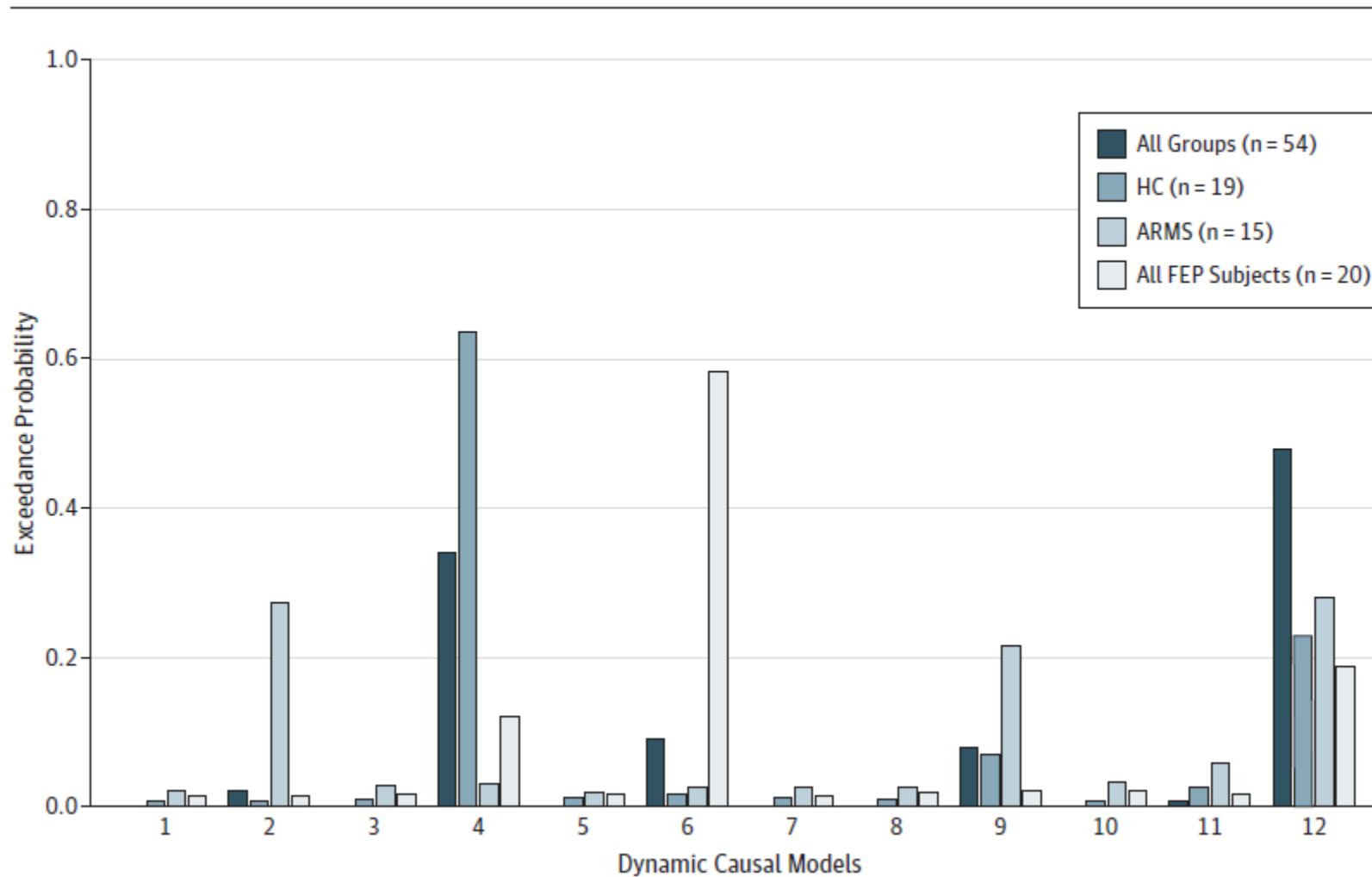


Prefrontal-parietal connectivity during working memory in schizophrenia

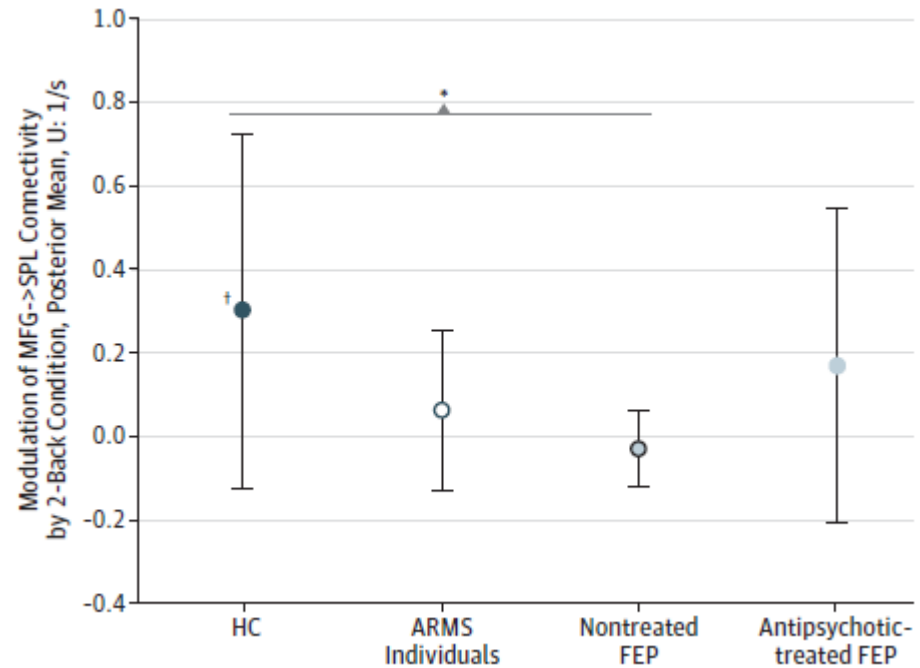
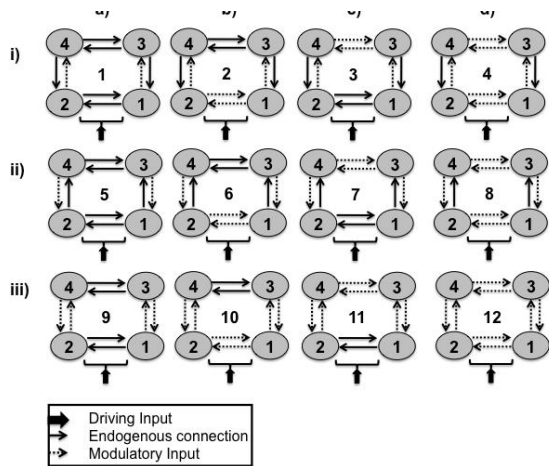
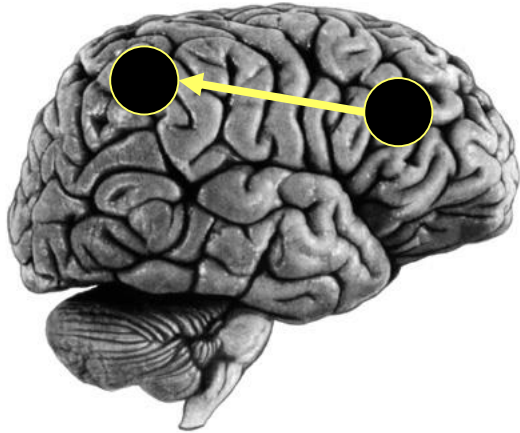


- 17 at-risk mental state (ARMS) individuals
- 21 first-episode patients (13 non-treated)
- 20 controls

BMS results for all groups



BMA results: PFC → PPC connectivity



17 ARMS, 21 first-episode (13 non-treated),
20 controls

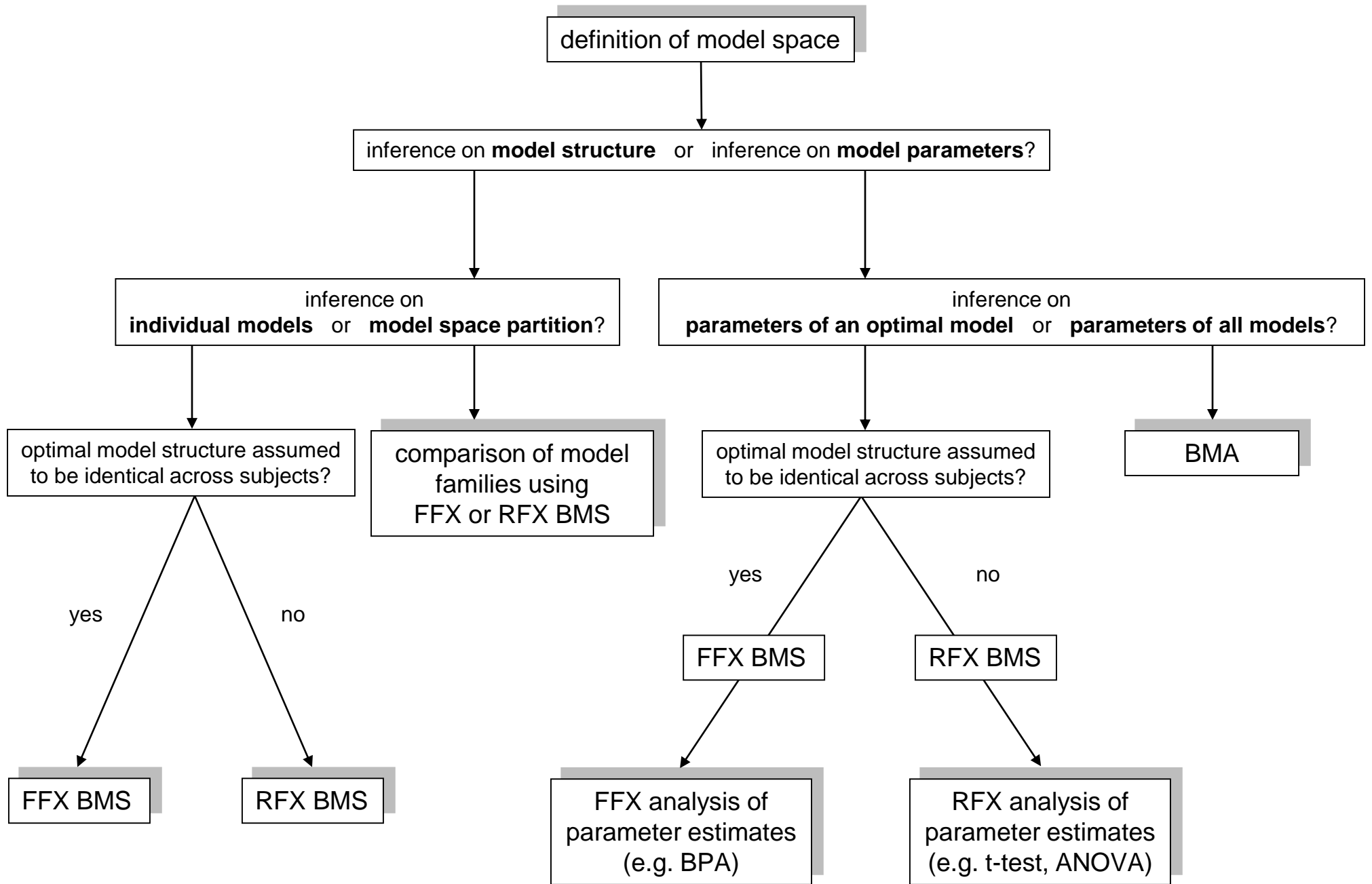
Protected exceedance probability: Using BMA to protect against chance findings

- EPs express our confidence that the posterior probabilities of models are different – under the hypothesis H_1 that models differ in probability: $r_k \neq 1/K$
- does not account for possibility "null hypothesis" H_0 : $r_k = 1/K$
- **Bayesian omnibus risk (BOR)** of wrongly accepting H_1 over H_0 :

$$P_0 = \frac{1}{1 + \frac{p(m|H_1)}{p(m|H_0)}}.$$

- **protected EP**: Bayesian model averaging over H_0 and H_1 :

$$\begin{aligned}\tilde{\varphi}_k &= P(r_k \geq r_{k' \neq k} | y) \\ &= P(r_k \geq r_{k' \neq k} | y, H_1)P(H_1 | y) + P(r_k \geq r_{k' \neq k} | y, H_0)P(H_0 | y) \\ &= \varphi_k(1 - P_0) + \frac{1}{K}P_0\end{aligned}$$



Further reading

- Penny WD, Stephan KE, Mechelli A, Friston KJ (2004) Comparing dynamic causal models. *NeuroImage* 22:1157-1172.
- Penny WD, Stephan KE, Daunizeau J, Joao M, Friston K, Schofield T, Leff AP (2010) Comparing Families of Dynamic Causal Models. *PLoS Computational Biology* 6: e1000709.
- Penny WD (2012) Comparing dynamic causal models using AIC, BIC and free energy. *Neuroimage* 59: 319-330.
- Rigoux L, Stephan KE, Friston KJ, Daunizeau J (2014) Bayesian model selection for group studies – revisited. *NeuroImage* 84: 971-985.
- Stephan KE, Weiskopf N, Drysdale PM, Robinson PA, Friston KJ (2007) Comparing hemodynamic models with DCM. *NeuroImage* 38:387-401.
- Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ (2009) Bayesian model selection for group studies. *NeuroImage* 46:1004-1017.
- Stephan KE, Penny WD, Moran RJ, den Ouden HEM, Daunizeau J, Friston KJ (2010) Ten simple rules for Dynamic Causal Modelling. *NeuroImage* 49: 3099-3109.
- Stephan KE, Iglesias S, Heinzle J, Diaconescu AO (2015) Translational Perspectives for Computational Neuroimaging. *Neuron* 87: 716-732.
- Stephan KE, Schlagenhauf F, Huys QJM, Raman S, Aponte EA, Brodersen KH, Rigoux L, Moran RJ, Daunizeau J, Dolan RJ, Friston KJ, Heinz A (2017) Computational Neuroimaging Strategies for Single Patient Predictions. *NeuroImage* 145: 180-199.

Thank you