# Bayesian inference and model inversion
## Variational + Markov chain Monte Carlo

## Lionel Rigoux

**Translational Neuro-Circuitry (TNC) Cologne**

**Translational Neuromodeling Unit (TNU) Zürich**

# Overview

Introduction to Bayesian inference (revisited)

Sampling methods (sampling)

Variational methods (approximation)

# Probability basics

Formalizes the degree of plausibility of events:

i.    represented by real numbers
ii.   should conform to intuition
iii.  should be consistent

Normalization

$$\sum_a p(a) = 1$$

Independence
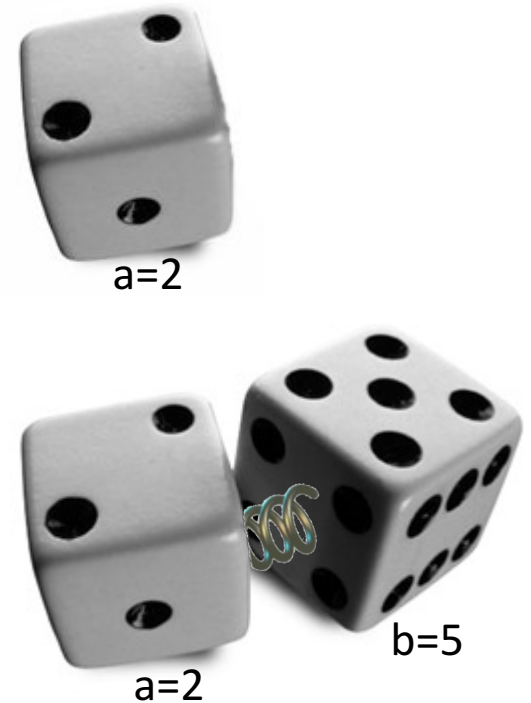
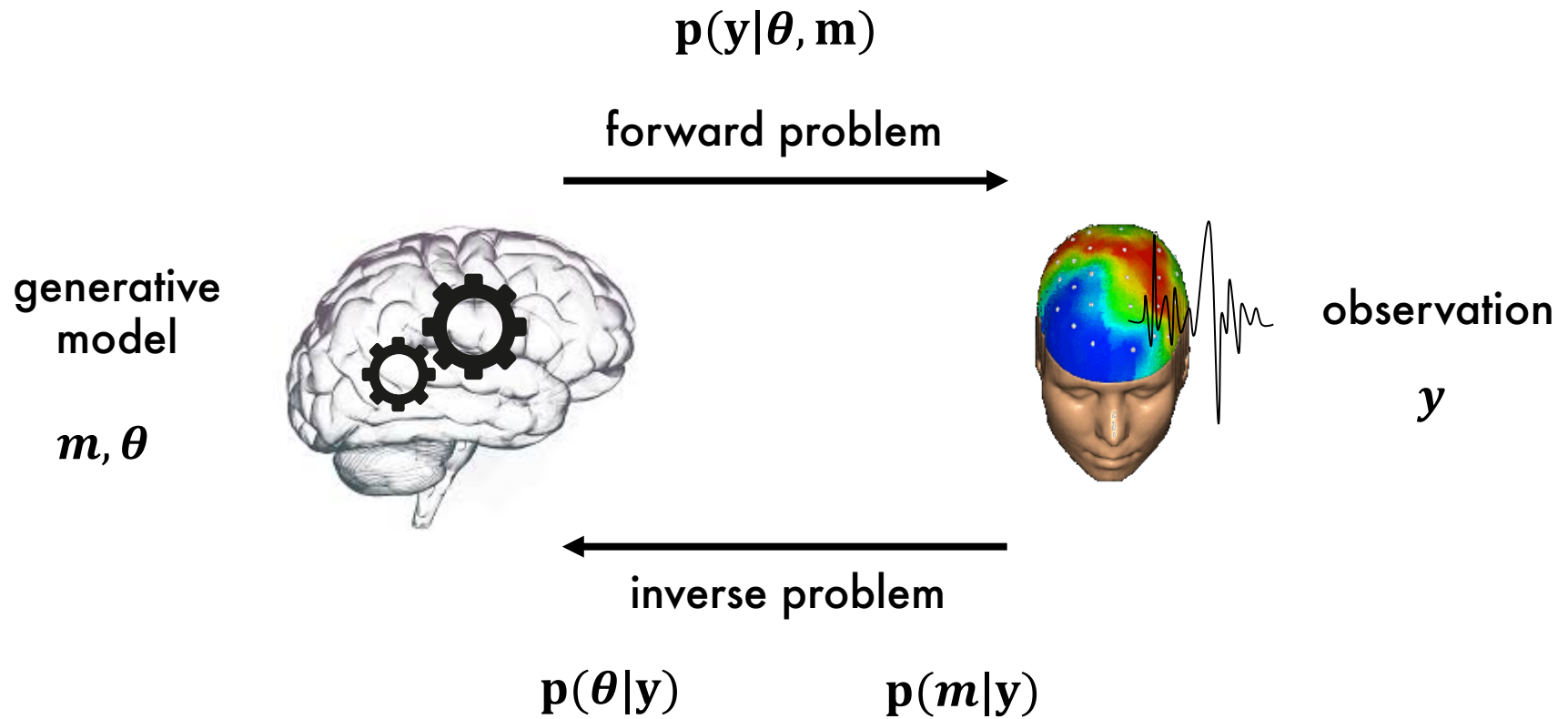$$\mathbf{p}(\mathbf{a}, \mathbf{b}) = \mathbf{p}(\mathbf{a})\mathbf{p}(\mathbf{b})$$

Conditioning

$$\mathbf{p}(\mathbf{a}, \mathbf{b}) = \mathbf{p}(\mathbf{a}|\mathbf{b})\mathbf{p}(\mathbf{b})$$

Marginalization

$$\mathbf{p}(\mathbf{b}) = \sum_a p(a, b)$$

a=2

b=5

a=2

# Model inversion

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

forward problem

generative
model

$$\boldsymbol{m}, \boldsymbol{\theta}$$

observation

$$\boldsymbol{y}$$

inverse problem

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

# Bayes rule

## Linear regression

| *model* | *prior* | *likelihood* | *posterior* |
|---|---|---|---|

$$\begin{aligned} y &= \theta\, x + \varepsilon \\ \varepsilon &= \mathcal{N}\left(0, \sigma_y^2\right) \end{aligned}$$

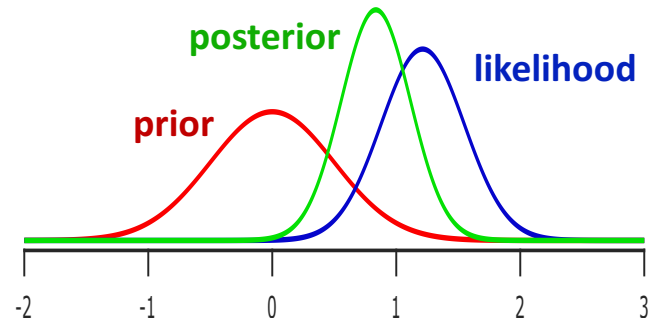$$p(\theta) = \mathcal{N}\left(0, \sigma_p^2\right)$$

$$p(y|\theta) = \mathcal{N}\left(\theta\, x, \sigma_y^2\right)$$
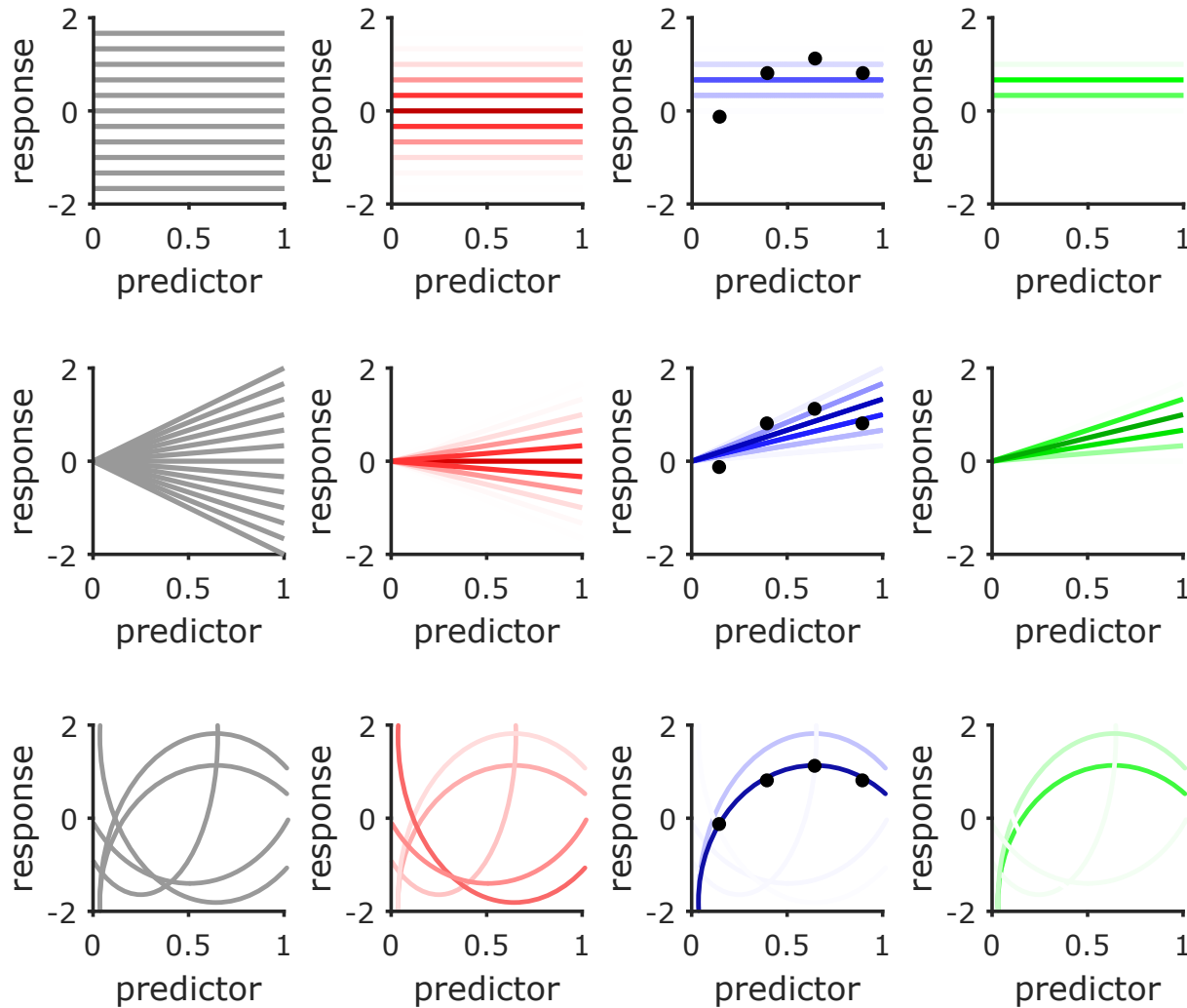
$$p(\theta|y)$$



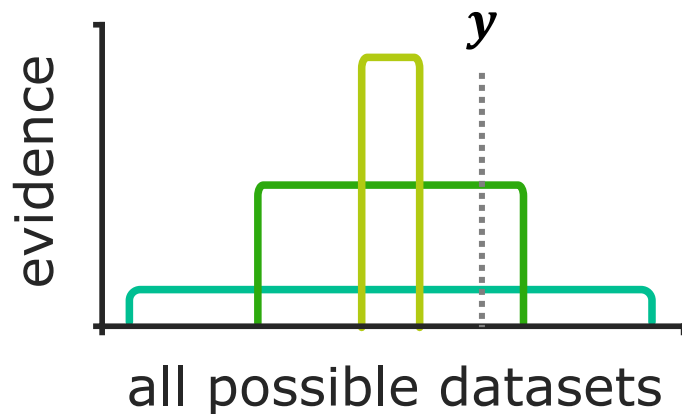$$p(\theta|y, m) = \frac{p(\theta|m)\, p(y|\theta, m)}{\int p(\theta|m) p(y|\theta, m)}$$



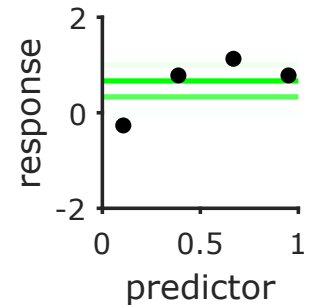posterior  likelihood  prior

# Model evidence

# Model evidence

Model evidence (marginal likelihood)
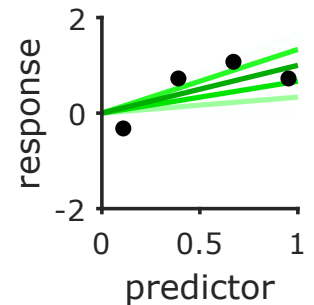
$$\int p(\theta|m)p(y|\theta, m)d\theta = p(y|m)$$

*"how likely are the data on average across plausible parametrization"*



too simple
miss the data

just right

too complex
overfitting

evidence

all possible datasets
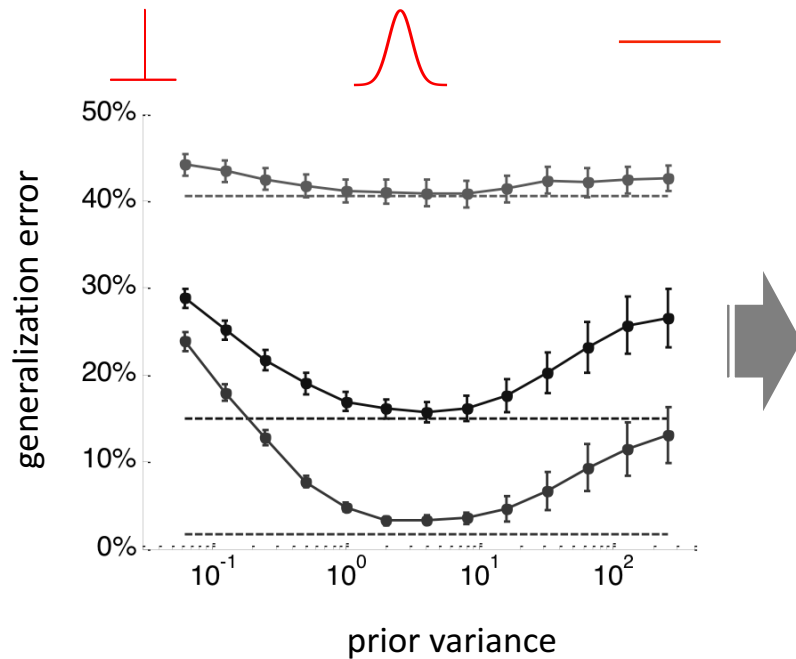
# The bias-variance trade-off



*Frequentist*

- always overfit (fit noise)
  > large variance
- estimate converge to the
  true value on average
  > unbiased

*Bayesian*

- regularized estimation
  > small variance
- estimate stays close to prior
  > biased

# Bayesian vs. frequentist hypothesis testing



parameter estimate

subject #

1-p = 1.000    ep = 1.000
t-test    Bayesian

ep = 1.000
1-p = 0.023
t-test    Bayesian

1-p = 0.999
ep = 0.500
t-test    Bayesian

1-p = 0.010
ep = 0.498
t-test    Bayesian

# On the importance of priors

Priors allow to define:
- plausible values of computational parameters
- range of data patterns predicted by the model

Role of priors
- avoid overfitting (generalization error)
- anchor a complexity measure

Impact of priors
- on parameters: "shrinkage to the mean" effect (bias / regularization)
- on model evidence

# Inference in practice

How to compute the posterior?

- Write $\textcolor{red}{p(\theta|m)}\,\textcolor{blue}{p(y|\theta,m)}$

1. recognize it looks like a know distribution

$$\textcolor{red}{p(\theta|m)} = \textcolor{red}{\mathcal{N}(\mu_0, \sigma_0^2)}$$
$$\textcolor{blue}{p(y|\theta,m)} = \textcolor{blue}{\mathcal{N}(\theta, \sigma_y^2)}$$

$\Rightarrow$

$$\textcolor{green}{p(\theta|y,m)} = \textcolor{green}{\mathcal{N}(\mu', \sigma^{2\prime})}$$
$$\textcolor{yellow}{p(y|m)} = \textcolor{yellow}{[\text{analytical solution}]}$$

$$\sigma^{2\prime} = \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma_y^2}\right)^{-1} \qquad \mu' = \left(\sigma^{2\prime}\right)^{-1}\left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum_n y_n}{\sigma_y^2}\right)$$

2. use variational Bayes or Monte Carlo methods

# Monte-Carlo methods

*If you can not calculate it,*
*simulate many random trials and see what happens...*

Law of large numbers

- simulate many indepedent draws of a random variable

- the average of the results will converge to the true expected value (probability mean)

$$E[\theta] = \int p(\theta|y, m)\theta \approx \frac{1}{n}\sum_n \theta_n$$

# A little game

The un-normalized posterior:

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

- is not a probability

- gives the relative plausibility of parameter values

# Markov Chain sampling
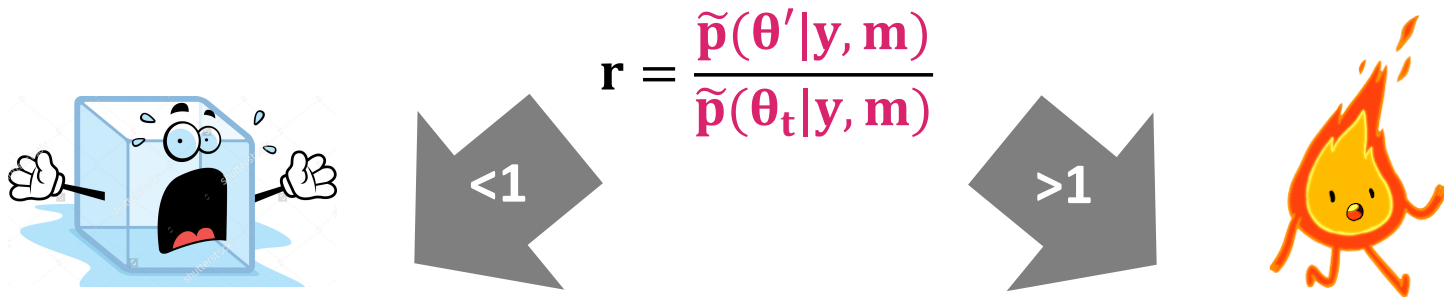
Markov Chain: stochastic process that evolve in time

- initial state $\theta_0$

- state evolve following a transition function $\mathrm{T}(\theta_{t+1}|\theta_t)$

> In the long run the probability of visiting $\theta$ is called the *ergodic density*

# Metropolis Hastings algorithm

The Metropolis-Hastings algorithm

- start form $\theta_0$

- propose a new value according to $T'(\theta'|\theta_t)$

- look for guidance

$$r = \frac{\widetilde{p}(\theta'|y, m)}{\widetilde{p}(\theta_t|y, m)}$$

**<1**

**>1**

$$\theta_{t+1} = \theta'$$

jump to proposed value

if $r > X \sim U(0, 1)$

$$\theta_{t+1} = \theta'$$

else

$$\theta_{t+1} = \theta_t$$

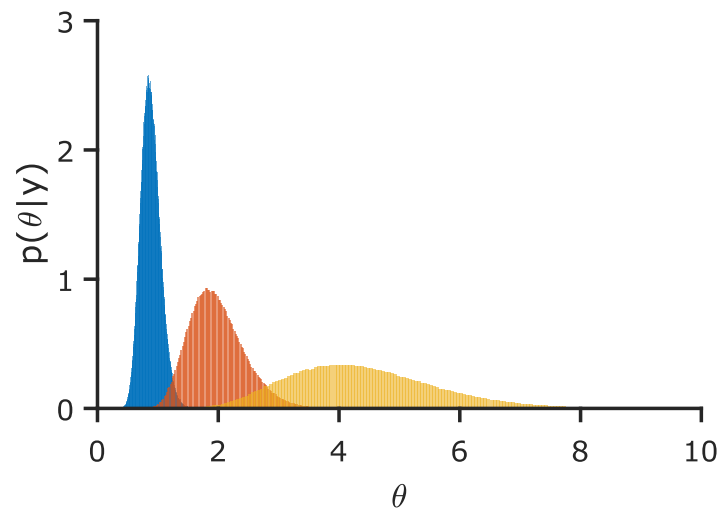ergodic density = $p(\theta|y, m)$

# Metropolis Hastings algorithm: example

Logistic regression

$$p(y = 1|\theta) = s(\theta x)$$

$$p(\theta) = \mathcal{N}(0, \sigma_0^2)$$

$$s = \frac{1}{1 + e^{-x}}$$

$$\tilde{p}(\theta|y) = \exp\left(-\frac{\theta^2}{2\sigma_0^2}\right) \times \prod_y s(\theta x)^y \left(1 - s(\theta x)\right)^{1-y}$$

# Metropolis Hastings algorithm: multivariate case

Example of linear regression

$$\mathbf{y} = \boldsymbol{\theta}\mathbf{x} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\lambda})$$



$$\mathbf{p}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbf{p}(\boldsymbol{\theta})\mathbf{p}(\boldsymbol{\lambda}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\sigma_0})\mathbf{Ga}(\mathbf{a}, \mathbf{b})$$

Blocked sampler

- start with $\theta_t = \theta_0$ and $\lambda_t = \lambda_0$

- repeat:

    - sample $\theta_{t+1}$ from $p(\theta|y, \lambda_t)$

    - sample $\lambda_{t+1}$ from $p(\lambda|y, \theta_{t+1})$

    Analytical expression (conjugacy)
    > no need for Markov chain!

- estimate $\mathbf{p}(\boldsymbol{\theta}, \boldsymbol{\lambda}|\mathbf{y})$ using LLN

# Monte-Carlo inference

Sample in turn from all the conditional (Gibbs) or the un-normalized conditional (Markov chain/Metropolis-Hastings) posterior.

> Sufficient statistics converge to the true value.

Problems:

- computationally expensive

- does not scale well

- no direct measure of model evidence
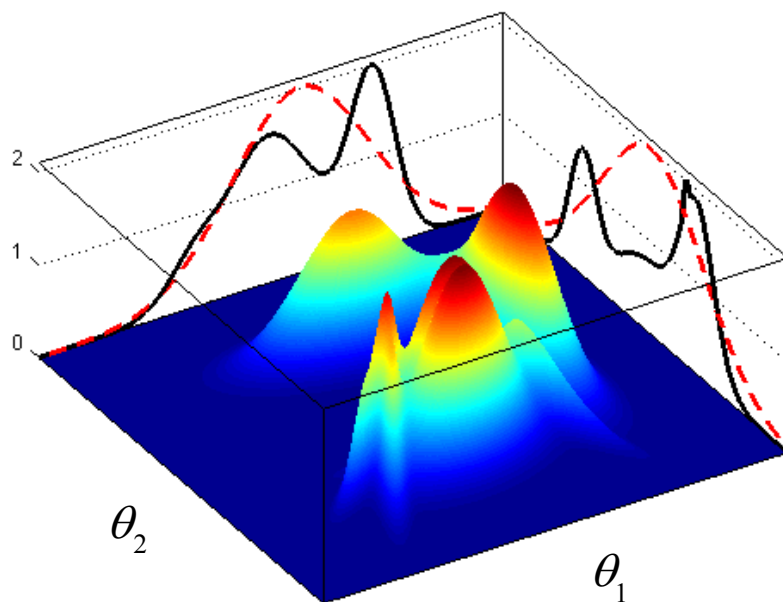
- hard to tune and diagnose

# Variational inference

*"variational inference is the thing you implement while you wait for your sampler to converge"*

David Blei

# Approximating the posterior

$$p(\theta_1, \theta_2 | y, m) = \frac{p(\theta_1, \theta_2 | m)\, p(y | \theta_1, \theta_2, m)}{p(y | m)}$$



Mean field approximation —

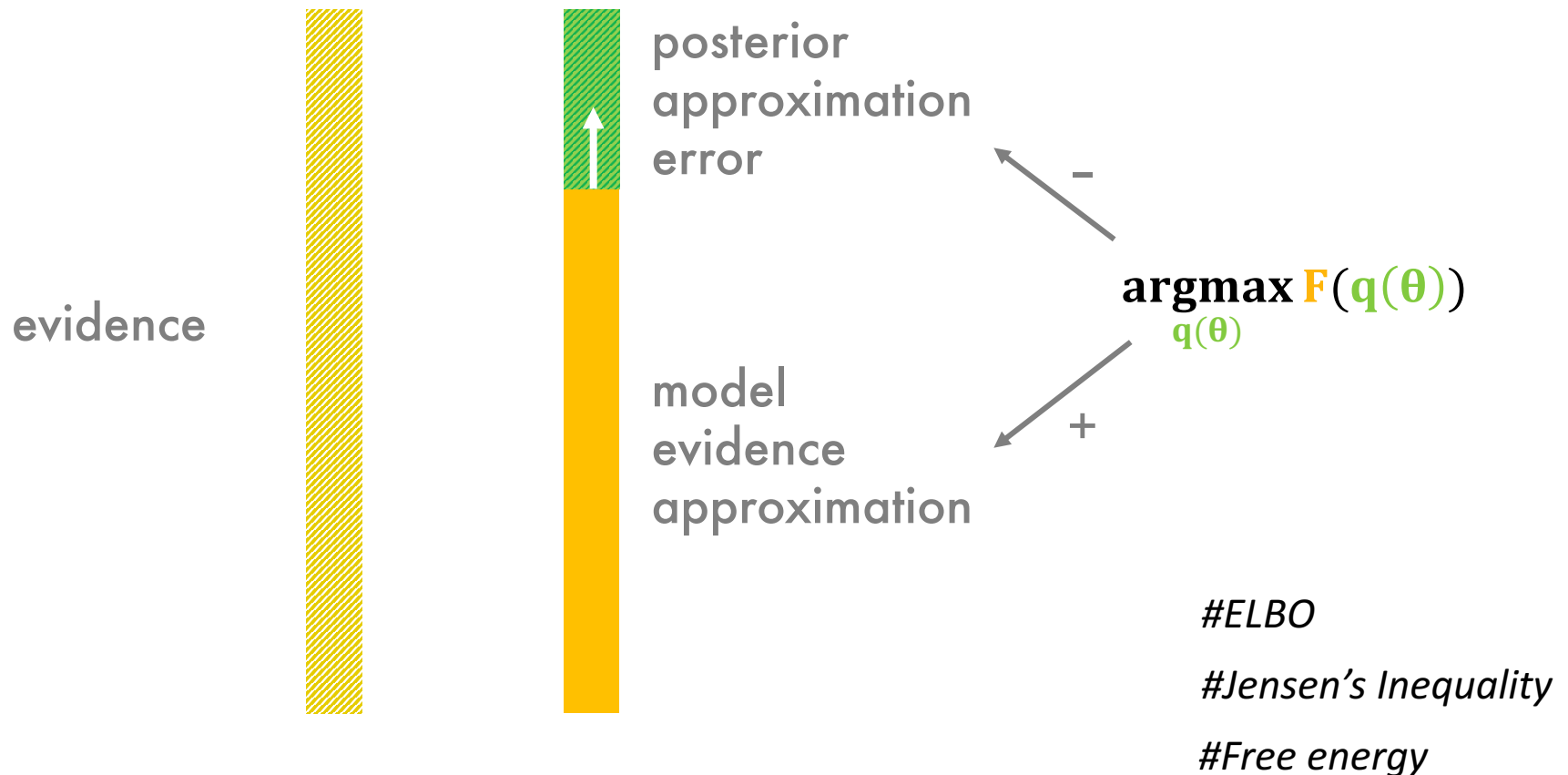$$p(\theta_1, \theta_2 | y) \approx p(\theta_1 | y) p(\theta_2 | y)$$

Laplace approximation - - -

$$q(\theta_1 | y) \approx \mathcal{N}(\mu_1, \Sigma_1)$$

finding $p(\theta_1, \theta_2 | y, m)$ ⟶ finding $\mu_1, \mu_2, \Sigma_1, \Sigma_2$

# Free Energy approximation

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

posterior
approximation
error

evidence

model
evidence
approximation

$$\mathbf{\underset{q(\theta)}{argmax} \, F(q(\theta))}$$

−

+

*#ELBO*

*#Jensen's Inequality*

*#Free energy*

# Free Energy approximation

Approximating the model evidence = maximizing the ELBO wrt $q(\theta)$

1) Maximize the free energy using variational calculus

$$F = \langle \log p(y|\theta_1, \theta_2, m) + \log p(\theta_1, \theta_2) \rangle_q + \langle \log q(\theta_1, \theta_2) \rangle_q$$

$$\frac{\partial F}{\partial q(\theta_1)} = 0 \implies q(\theta_1) \propto \exp\left(\langle \log p(y|\theta_1, \theta_2, m) + \log p(\theta_1, \theta_2) \rangle_{q(\theta_2)}\right)$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{variational free energy} \quad I(\theta_1)}$$

2) Iterate over parameters until convergence

# Variational inference

Find $\mathcal{N}(\mu, \Sigma)$ that best approximate $I(\theta)$

## logistic regression



$$I(\theta) = \exp\left(-\frac{\theta^2}{2\sigma_0^2}\right) \times \prod_y s(\theta x)^y \left(1 - s(\theta x)\right)^{1-y}$$

## multivariate case

Until convergence:

for all i:

➢ $\mu_i = \max_{\theta_i}\left(I(\theta_i)\right)$

➢ $\Sigma_i = -\left[\left.\frac{\partial^2}{\partial\theta_i^2}\right|_{\mu_i} I(\theta_i)\right]^{-1}$

end
end

# Variational inference

Summarize the posterior to its sufficient statistics (mean, variance) and optimize those values wrt the evidence lower bound.

This requires multiple approximations (free-energy, mean-field, Laplace) to be tractable.

Problems:

- does not converge to the true posterior
- can get stuck in local optimum

Model evidence (normalization factor of the posterior) is in general intractable.

Sampling methods give a computationally expensive estimation of the true posterior.

Variational methods are fast and scalable but potentially inaccurate.

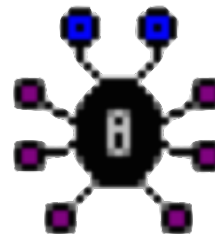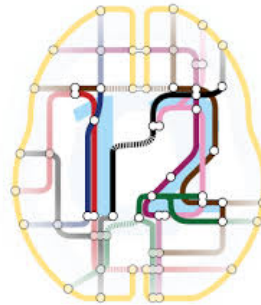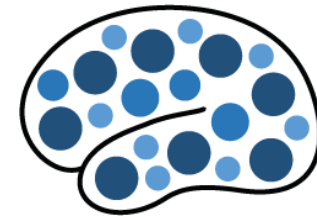# Software

Variational
  VBA-toolbox
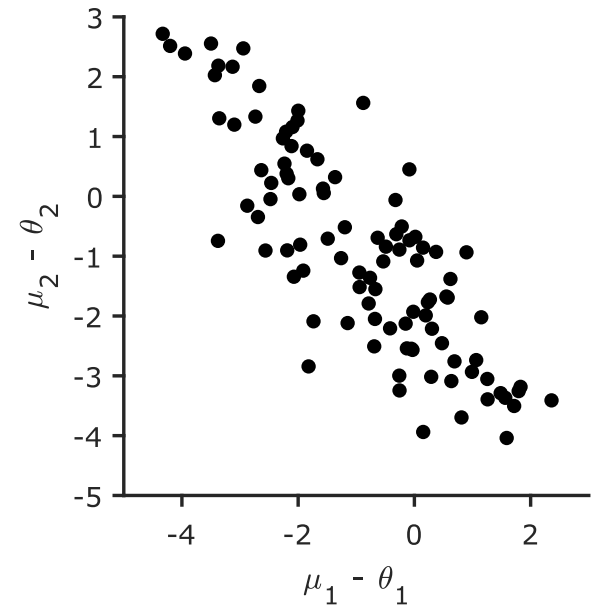  TAPAS
  SPM
Sampling
  STAN
  BUGS
  JAGS
  hBayesDM
  hddm
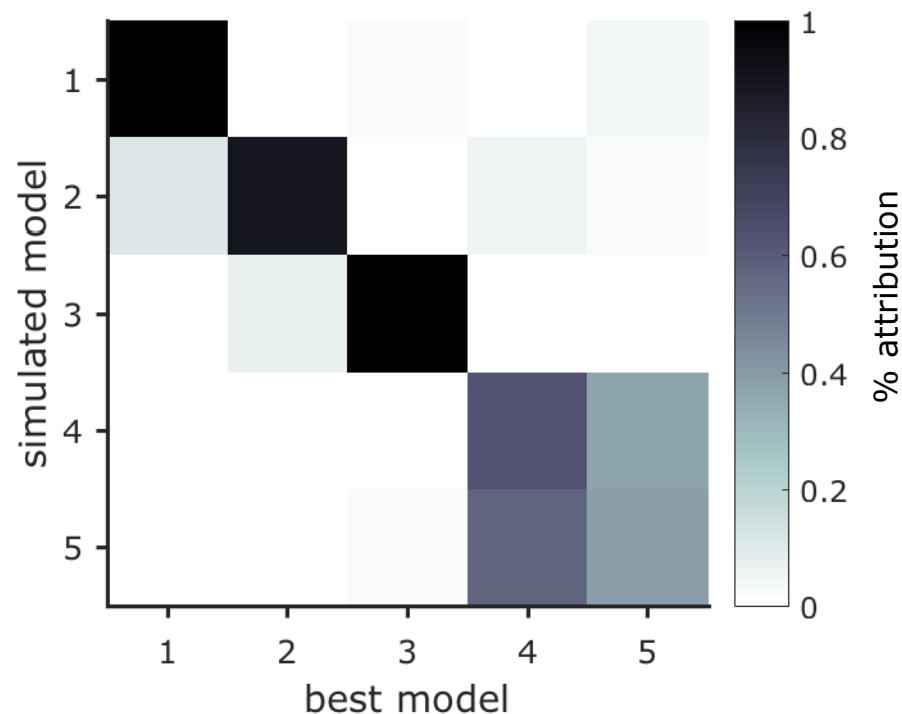
JAGS

Checking if your parameters are identifiable:

- simulate data using your

  design with realistic $\theta$

- invert your model (find $\mu$)

- compute estimation error ($\mu - \theta$)

- check bias/variance trade-off

- check for posterior / error correlation

Checking if your models are identifiable:

- simulate all models

- compute evidence of each hypothesis for each dataset (BMS)

- count misattributions and build confusion matrix

# Thank you!