# Fitting a model: VB & MCMC
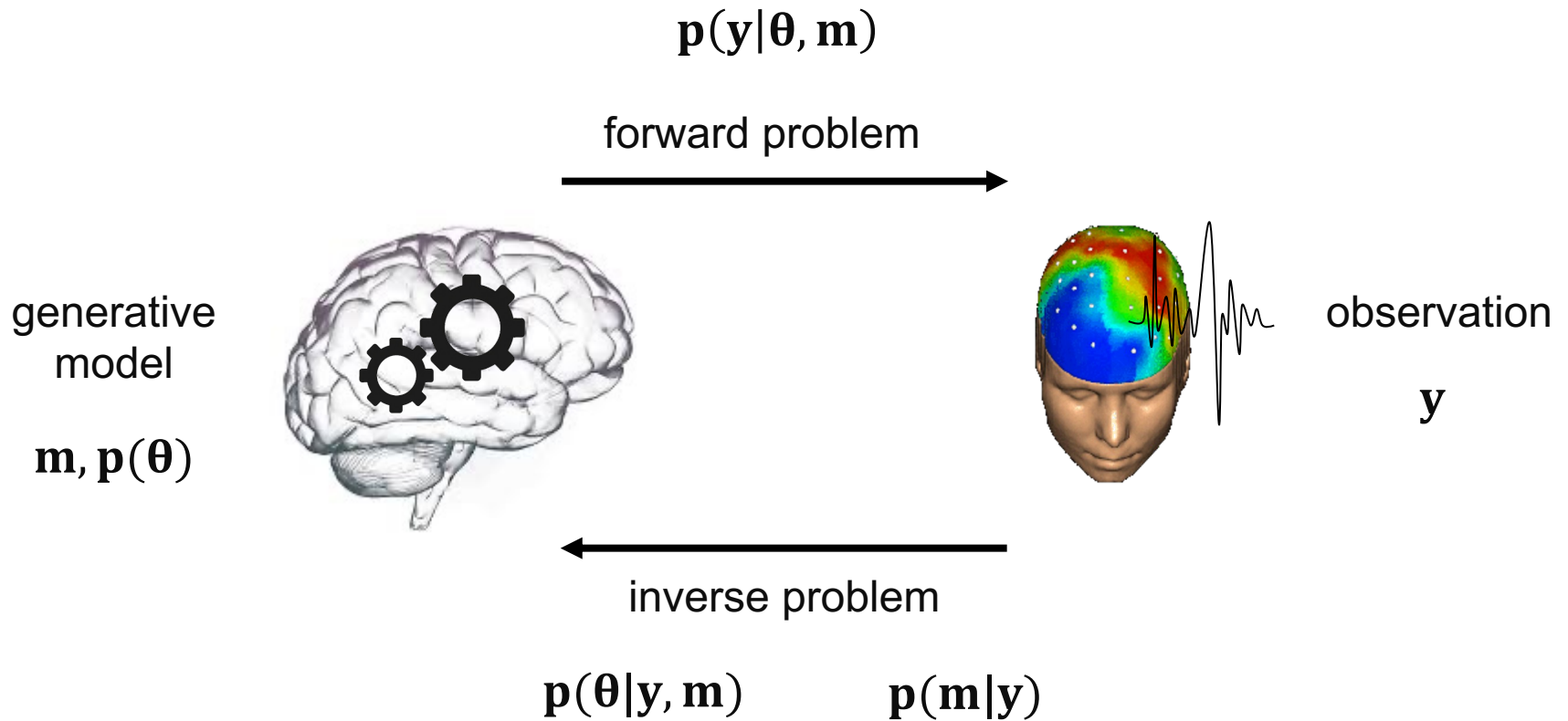
## Lionel Rigoux

Max Planck Institute for Metabolism Research

Translational Neuro-Circuitry Group

Max-Planck-Institut
für Stoffwechselforschung

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

forward problem

generative
model

$$m, p(\theta)$$

observation

$$y$$

inverse problem

$$p(\theta|y, m) \qquad p(m|y)$$

# Bayes rule

**Joint distribution**

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

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

**Expectation**

$$E[p(y | \theta, m)]_{p(\theta | m)}$$

**Marginal likelihood**

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

**Model evidence**

$$p(y | m)$$

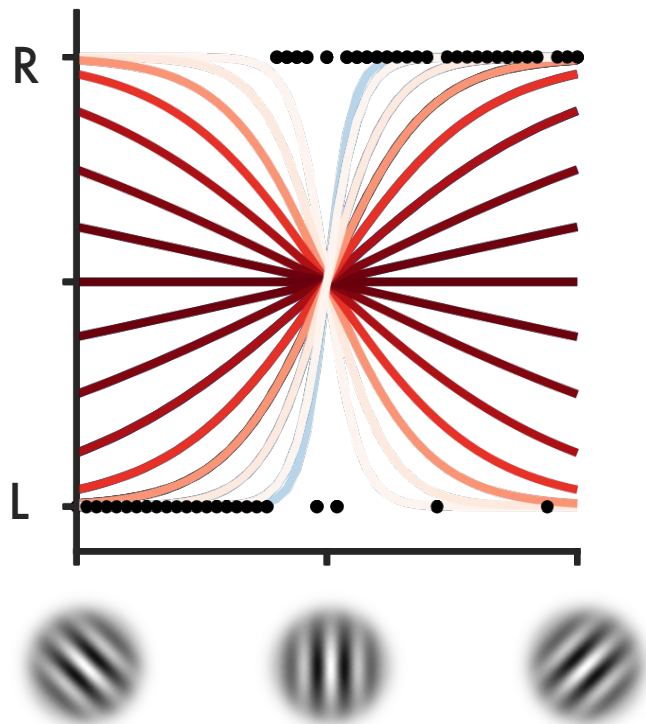Compute the posterior and the evidence for a model

Monte-Carlo (sampling) methods

Variational methods

Good practices

github.com/lionel-rigoux/tutorial-bayesian-inference

# Example: logistic regression



*Sensitivity to orientation?*

*Bias?*

**Model prediction**

$$p(y = 1|u, \theta, \beta) = sig(\theta u + \beta) = s$$

**Likelihood**

$$\log p(y|\theta, \beta) = \sum y \log s + (1 - y) \log(1 - s)$$

**Prior**

$$\theta \sim \mathcal{N}(\mu_\theta, \sigma_\theta^2) \qquad \beta \sim \mathcal{N}(0, 0)$$

$$\log p(\theta) = -\frac{1}{2} \left[ \frac{(\theta - \mu_\theta)^2}{\sigma_\theta^2} + \log 2\pi\sigma^2 \right]$$

**Joint**

$$\log p(y, \theta, \beta) =$$

$$\sum y \log s + (1 - y) \log(1 - s) - \frac{(\theta - \mu_\theta)^2}{2\,\sigma_\theta^2} + cst$$

# Example: logistic regression

**Joint**

$$\log \mathbf{p}(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \sum y \log s + (1 - y) \log(1 - s) - \frac{(\boldsymbol{\theta} - \boldsymbol{\mu}_{\boldsymbol{\theta}})^2}{2\,\sigma_{\boldsymbol{\theta}}^2} + cst$$

$$\mathbf{p}(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) \propto \prod s^{\,y}\,(1 - s)^{\,1-y}\,e^{-\frac{(\boldsymbol{\theta} - \boldsymbol{\mu}_{\boldsymbol{\theta}})^2}{2\,\sigma_{\boldsymbol{\theta}}^2}}$$

**Posterior**

$$\mathbf{p}(\boldsymbol{\theta}, \boldsymbol{\beta} | \mathbf{y}) \propto \mathbf{p}(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) \qquad\qquad \mathbf{MAP} = \operatorname*{argmax}_{\theta, \beta} \mathbf{p}(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta})$$

**Model evidence**

$$\mathbf{p}(\mathbf{y}) = \int \mathbf{p}(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta})\, d\boldsymbol{\theta} d\boldsymbol{\beta}$$

# Sampling (Monte Carlo)

# Monte-Carlo methods

**Expectation (theoretical mean)**

$$E[z] = \sum p(z)z = \sum_{z=1}^{6} \frac{1}{6}z = 3.5$$

**Variance (theoretical distance to the mean)**

$$E[(z - 3.5)^2] = \sum p(z)(z - 3.5)^2 = 2.9167$$

**Expectation $\approx$ Empirical mean**

$$E[z] \approx \frac{1}{n} \sum_{i=1}^{n} z_i \qquad z_i \sim p(z)$$

$$E[f(z)] \approx \frac{1}{n} \sum_{i=1}^{n} f(z_i)$$

Law of
Large Numbers

# Monte-Carlo methods

**Model evidence**

*Arithmetic estimator*

$$\mathbf{p}(\mathbf{y}) = \mathbf{E}[\mathbf{p}(\mathbf{y}|\boldsymbol{\theta})]_{\mathbf{p}(\boldsymbol{\theta})} \approx \frac{1}{\mathbf{n}} \sum \mathbf{p}(\mathbf{y}|\boldsymbol{\theta}_i)$$

*Samples from prior*

$$\boldsymbol{\theta}_i \sim \mathbf{p}(\boldsymbol{\theta})$$

*Harmomic estimator, Gibb's estimator,*
*Annealed importance sampling, etc.*

**Posterior moments**

*Mean*

$$\boldsymbol{\mu} = \mathbf{E}[\boldsymbol{\theta}]_{\mathbf{p}(\boldsymbol{\theta}|\mathbf{y})} \approx \frac{1}{\mathbf{n}} \sum \boldsymbol{\theta}_i$$

*Samples from posterior*

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

*Variance*

$$\boldsymbol{\Sigma} = \mathbf{E}\left[(\boldsymbol{\theta} - \boldsymbol{\mu})^2\right]_{\mathbf{p}(\boldsymbol{\theta}|\mathbf{y})} \approx \frac{1}{\mathbf{n}} \sum (\boldsymbol{\theta}_i - \widehat{\boldsymbol{\mu}})^2$$
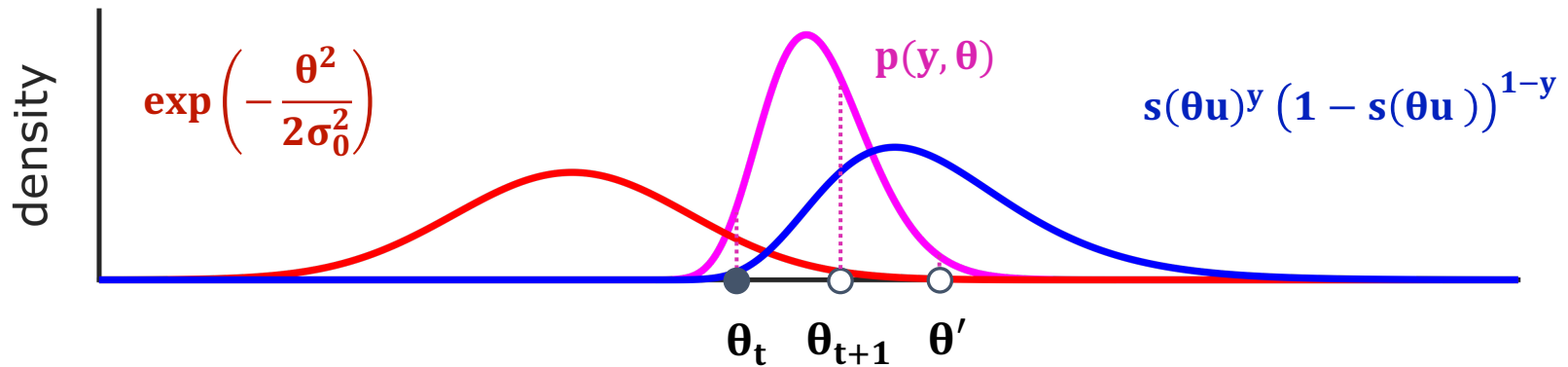
# A little game

The joint as an un-normalized posterior:

$$p(\theta|y) \propto p(\theta)\, p(y|\theta) = p(\theta, y)$$

- is not a probability over parameters
- gives the relative plausibility of parameter values

# Metropolis-Hastings algorithm



density

$$\exp\left(-\frac{\theta^2}{2\sigma_0^2}\right)$$

$p(y, \theta)$

$$s(\theta u)^y \left(1 - s(\theta u)\right)^{1-y}$$

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

**Current state**

$$p(y, \theta_t) = p(\theta_t)\, p(y|\theta_t)$$

**Proposal**

$$\theta' \sim q(\theta|\theta_t)$$

$$p(y, \theta') = p(\theta')\, p(y|\theta')$$

$$\alpha = \frac{p(y, \theta')}{p(y, \theta_t)}$$

$\alpha \geq 1$

*Jump to proposed value*

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

*Draw* $x \sim U(0, 1)$

- if $\alpha > x$, jump
  $$\theta_{t+1} = \theta'$$
- else, stay in place
  $$\theta_{t+1} = \theta_t$$

$\alpha < 1$

# Did I sample right?

All sampling methods requires some "post-processing" and an extensive diagnostic to ensure the samples are representative.



1) Run multiple chains

2) Check:

- Convergence (eg. Geweke)

- Mixing (eg. Gelman-Rubin)

- Autocorrelation (decimation)

- Step size (Goldilocks principle)

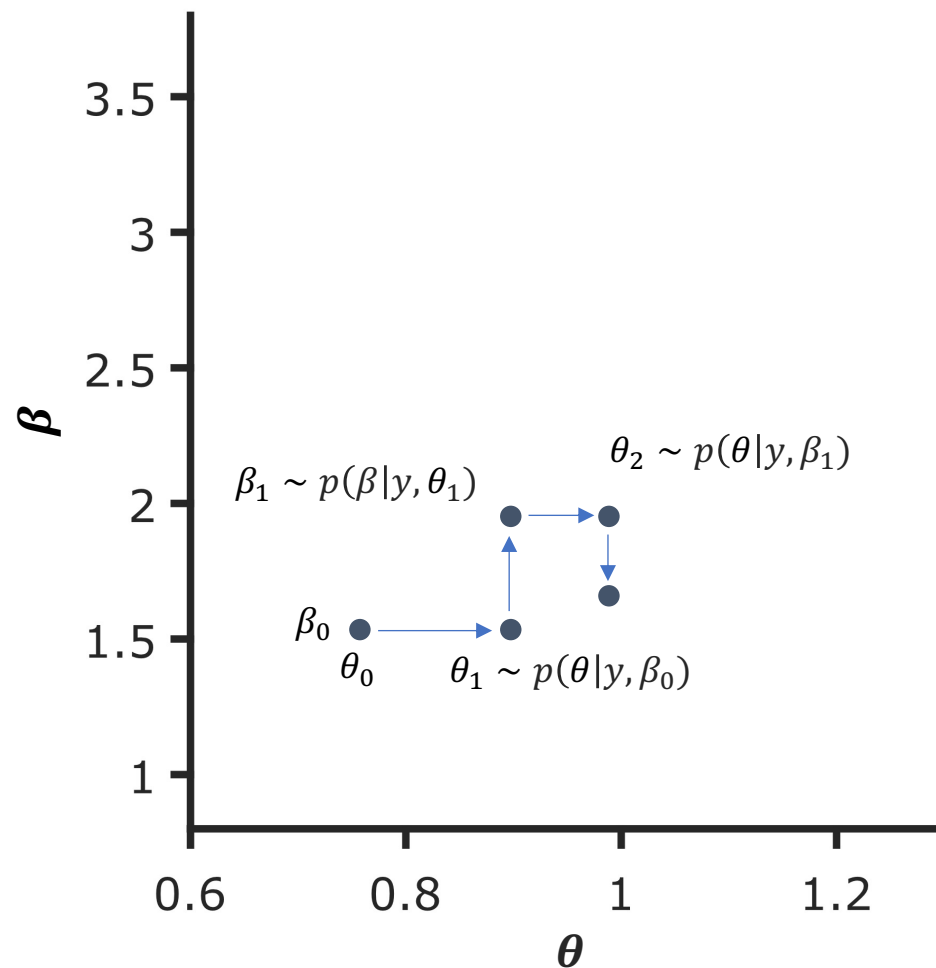# Multivariate case

**Write conditional posteriors**

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

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

**Iterative sampling**

$$\theta_t \sim p(\theta|y, \beta_{t-1})$$

$$\beta_t \sim p(\beta|y, \theta_t)$$

# Multivariate case

Using the law of large numbers:

**Posterior mean**

$$E[\theta|y] \approx \text{mean}(\theta_t)$$

$$E[\beta|y] \approx \text{mean}(\beta_t)$$

**Posterior variance**

$$E\big[(\theta - \overline{\theta})^2\big|y\big] \approx \text{var}(\theta_t)$$

$$E\left[(\beta - \overline{\beta})^2\big|y\right] \approx \text{var}(\beta_t)$$

**Covariance, etc.**

# Monte-Carlo inference

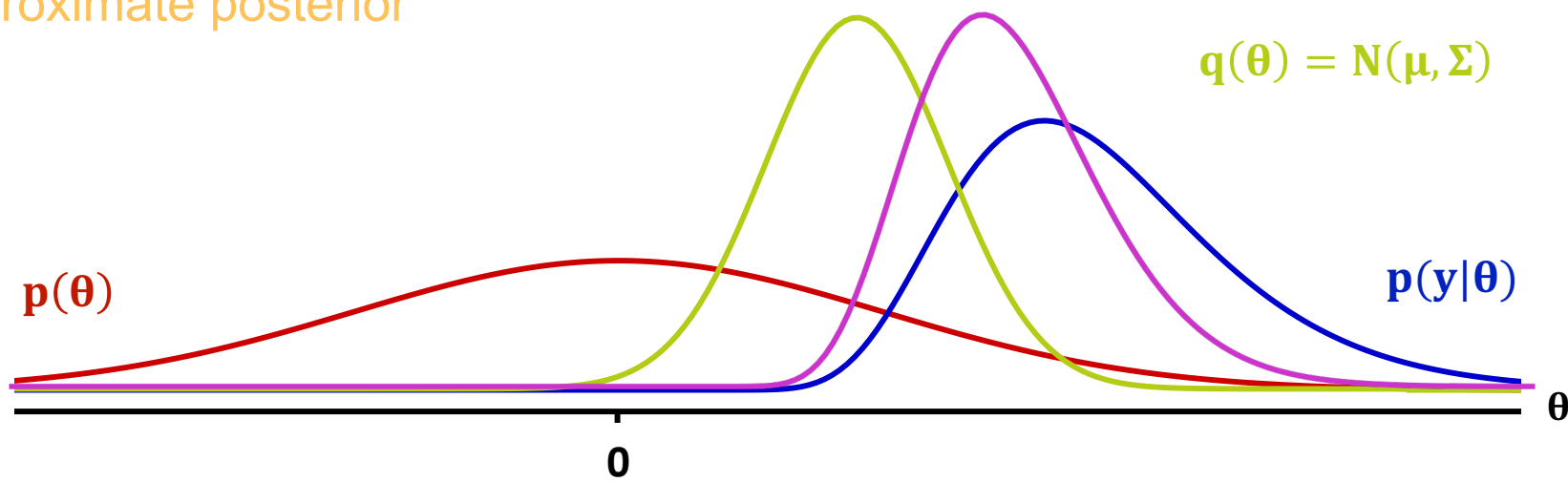Monte-Carlo methods rely on sampling to estimate the posterior and the model evidence.

The Law of Large Numbers guarantees that the sufficient statistics of the samples will converge to the true posterior moments.
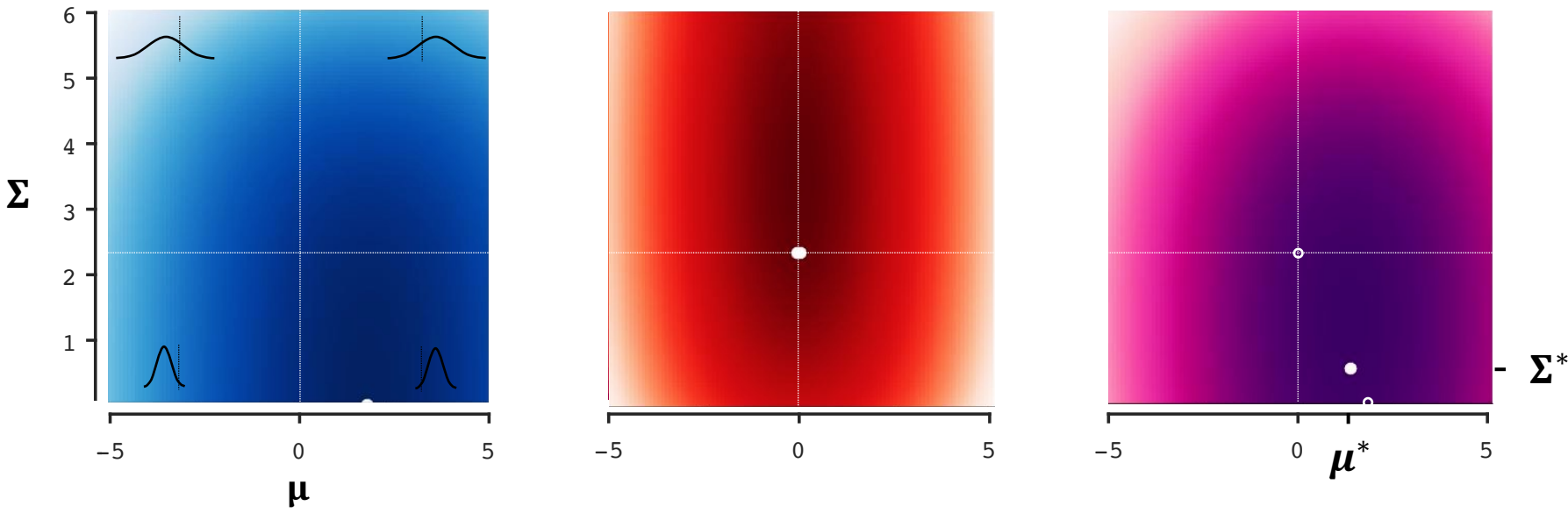
**Problems**

- computationally expensive
- does not scale well with the number of parameters
- hard to tune and diagnose
- no direct measure of model evidence

# Variational Methods

# Approximate posterior

$q(\theta) = N(\mu, \Sigma)$

$p(\theta)$

$p(y|\theta)$

$\theta$

$0$

$$\mathbf{E}[\log \mathbf{p(y|\theta)}]_{\mathbf{q}} \quad + \quad \mathbf{E}\left[\log \frac{\mathbf{p(\theta)}}{\mathbf{q(\theta)}}\right]_{\mathbf{q}} \quad = \quad \mathbf{E}\left[\log \frac{\mathbf{p(y,\theta)}}{\mathbf{q(\theta)}}\right]_{\mathbf{q}}$$

$\Sigma$

$\mu$

$- \Sigma^*$

$\mu^*$

**candidate distribution**     $q(\theta)$

**Jensen's inequality**

$$\log p(y) = \log \int p(y, \theta) \, d\theta$$

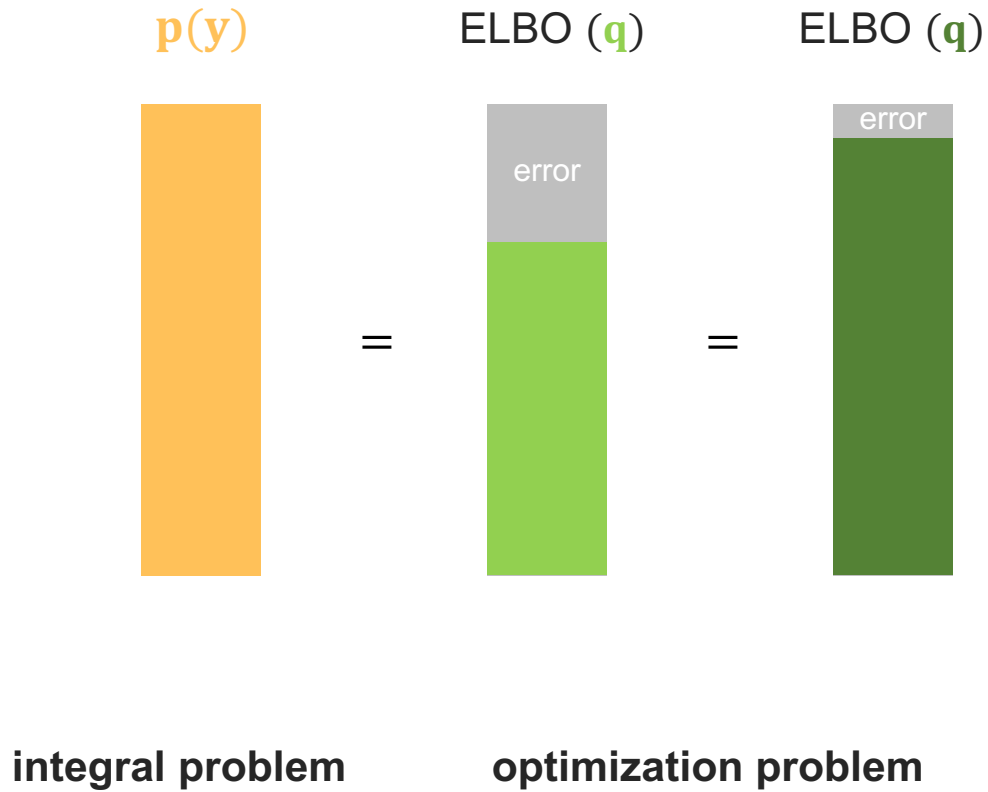$$= \log \int \frac{p(y, \theta)}{q(\theta)} q(\theta) \, d\theta$$

$$= \log E \left[ \frac{p(y, \theta)}{q(\theta)} \right]_{q(\theta)}$$

$$= E \left[ \log \frac{p(y, \theta)}{q(\theta)} \right]_{q(\theta)} + KL[q(\theta) || p(\theta | y)]$$

$$\text{ELBO} \qquad\qquad\qquad \text{error}$$

$$< p(y) \qquad\qquad\qquad > 0$$

# Evidence LOwer Bound



$p(y)$    ELBO ($q$)    ELBO ($q$)

**integral problem**    **optimization problem**

# Maximizing the ELBO

$$\log \mathbf{p(y)}$$
$$\approx \max \mathbf{E}\left[\log \frac{\mathbf{p(y,\theta)}}{\mathbf{q(\theta)}}\right]_{\mathbf{q(\theta)}}$$

**Variational Laplace**

*Using exponental family*

$$\mathbf{q(\theta)} = N(\mathbf{\mu}, \mathbf{\Sigma})$$

*Analytical approximation*

$$\mathbf{ELBO} \approx \mathbf{ELBO_{Laplace}}$$

*Find maximum*

$$\frac{d}{d\mathbf{q(\theta)}}\mathbf{ELBO_{Laplace}} = \mathbf{0}$$

**Solution**

$$\mathbf{\mu}^* = \operatorname{argmax} \mathbf{p(y, \theta)} \text{ = MAP}$$

$$\mathbf{\Sigma}^* = -\left[\left.\frac{\partial^2}{\partial\mathbf{\theta}^2}\right|_{\mathbf{\mu}^*} \log \mathbf{p(y, \theta)}\right]^{-1}$$

$$\log \mathbf{p(y)} \approx \log \mathbf{p(y, \mu}^*) + \frac{1}{2}[\log |\mathbf{\Sigma}^*| + \mathbf{n_\theta}\log(2\mathbf{\pi})]$$

# Multivariate posterior

**Mean field approximation**
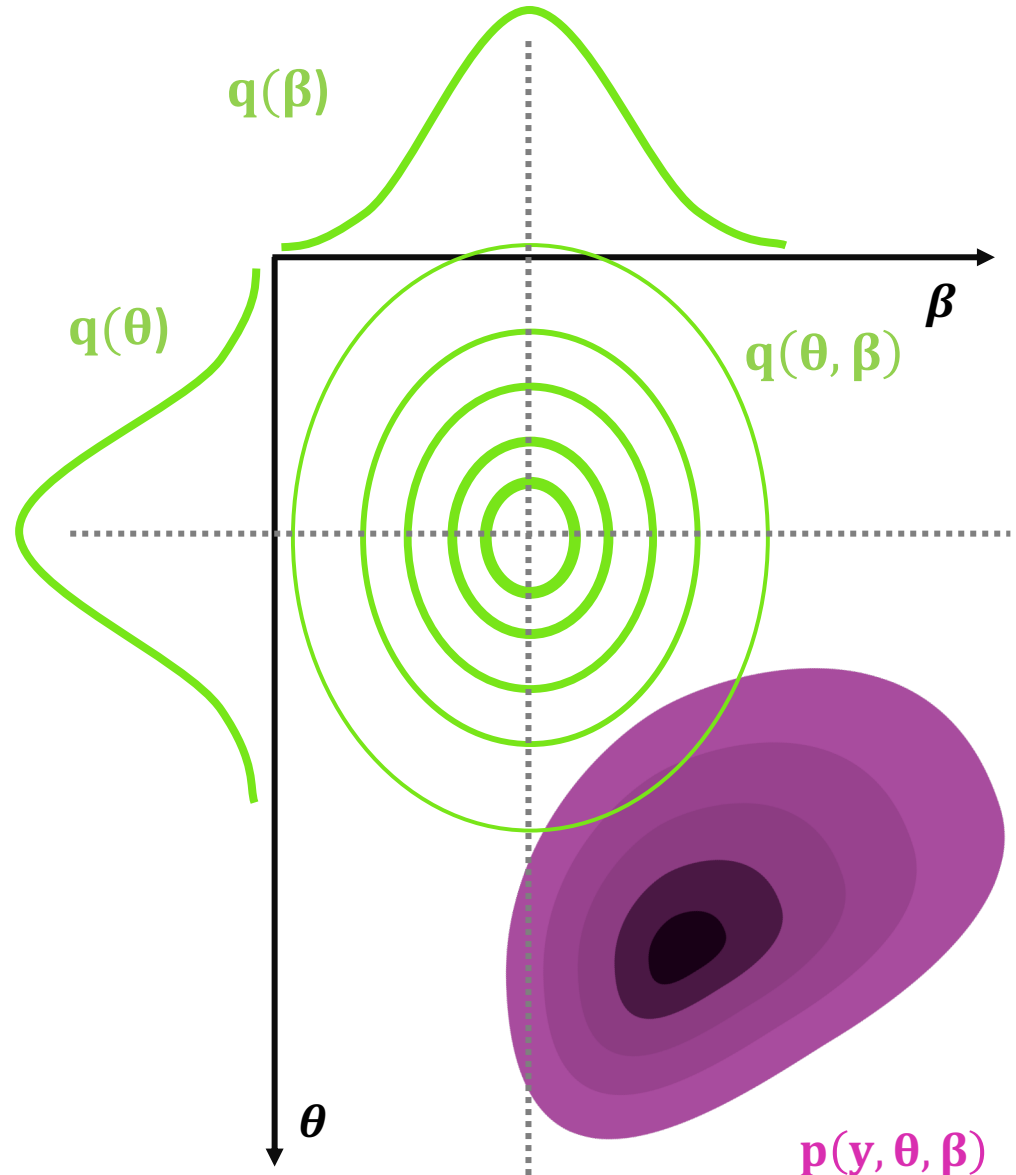
$$q(\theta, \beta) \approx q(\theta)q(\beta)$$

**Variational energy**

$$I(\theta) = E[\log p(y, \theta, \beta)]_{q(\beta)}$$

$$\approx \log p(y, \theta, \mu_\beta) + \ldots$$

**Iterative optimization**

$$\mu_i = \operatorname{argmax} I(\theta_i)$$

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

$q(\beta)$

$q(\theta)$

$q(\theta, \beta)$

$\beta$

$\theta$

$p(y, \theta, \beta)$

# Multivariate posterior

**Mean field approximation**

$$q(\theta, \varphi) \approx q(\theta)q(\beta)$$

**Maximise Variational energy**

$$I(\theta) = E[\log p(y, \theta, \beta)]_{q(\beta)}$$

$$\approx \log p(y, \theta, \mu_\beta) + \dots$$

**Iterative optimization**

$$\mu_i = \text{argmax } I(\theta_i)$$

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

$q(\beta)$

$\beta$

$q(\theta, \beta)$

$q(\theta)$

$\theta$

$p(y, \theta, \beta)$

$$\log \mathbf{p(y)}$$

$$\approx \max \mathbf{E}\left[\mathbf{log}\ \frac{\mathbf{p(y,\theta)}}{\mathbf{q(\theta)}}\right]_{\mathbf{q(\theta)}}$$

## Stochastic gradient

*Using samplable distribution*

$$\mathbf{q(\theta)} = N(\mathbf{\mu}, \mathbf{\Sigma})$$

*Gradient*

**∇ELBO**

$$= \mathbf{E}\left[\mathbf{\nabla}\log \mathbf{q(\theta)}\left(\log \frac{\mathbf{p(y,\theta)}}{\mathbf{q(\theta)}}\right)\right]_{\mathbf{q(\theta)}}$$

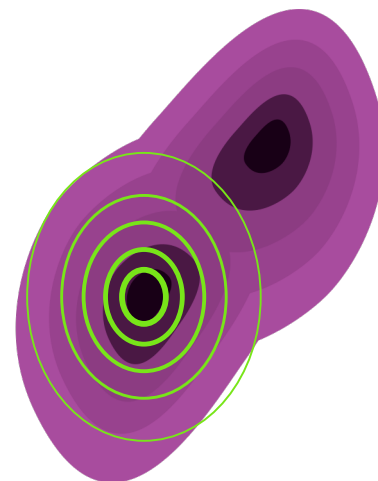**Solution**

*Ascend (MC approximation)*

# Variational inference

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

This requires multiple approximations (Jensen/Free-energy, Gaussian posterior, Laplace, mean-field) to be tractable.

**Problems**

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

# Take home message

Model evidence (normalization factor of the posterior) is in general intractable and calls for numerical methods.

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

✓ Variational methods are fast & scalable computations of an approximation of the posterior.
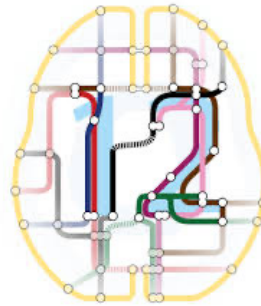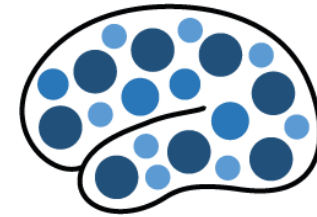
✓ Other techniques in development: Deep Bayesian Inversion

# Software

**Variational**
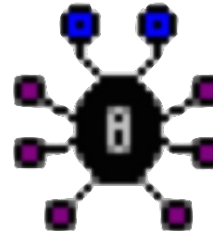
    VBA-toolbox

    TAPAS

    SPM

**Sampling**

    STAN

    BUGS

    JAGS

    hBayesDM

    hddm

# VBA Toolbox

282 published papers

85 demos (tutorial, Q-learning, HGF, DCMs, etc)

Online wiki + Q&A

Simulation

Inversion (single subject, hierarchical)

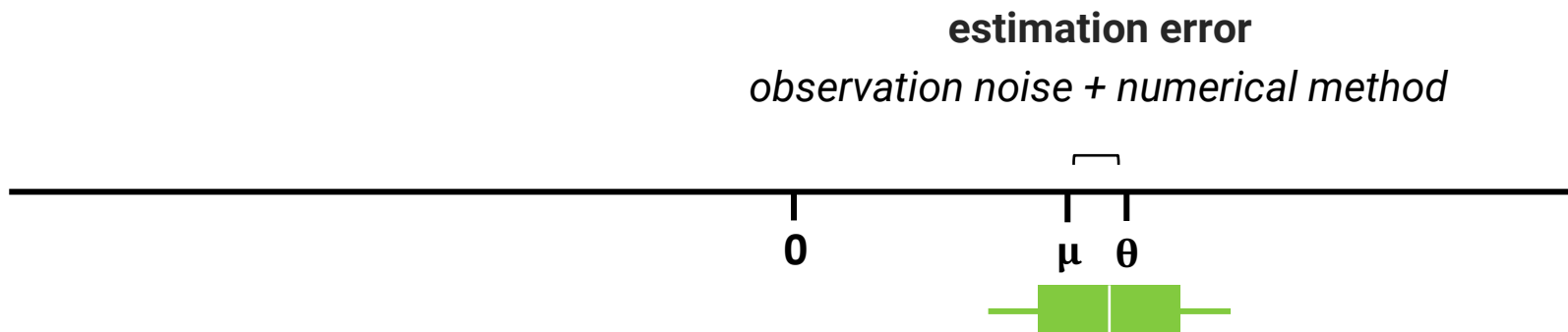Model selection (families, btw groups, btw conditions)

Visual diagnostics

Design optimization, multisession, multimodal observations, …

Need only the model description!

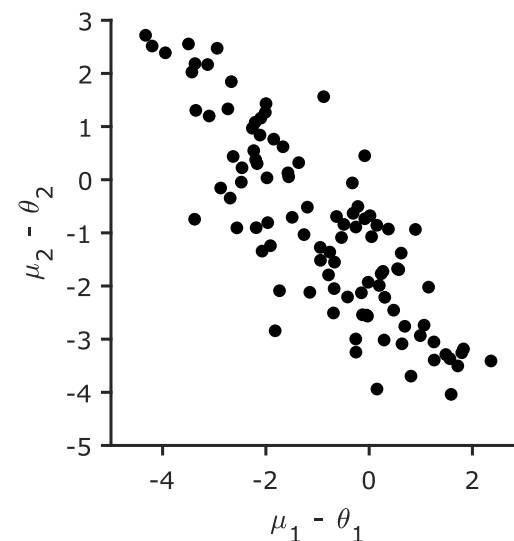# Validating your pipeline: parameter identifiability

**estimation error**
*observation noise + numerical method*



simulate data using your design with a realistic $\boldsymbol{\theta}$

do check if model predictions do emerge

invert your model (find $\boldsymbol{\mu}$)

compute estimation error ($\boldsymbol{\mu} - \boldsymbol{\theta}$)

- check effect of prior mean
- check effect of prior variance
- assess overfitting
- check for posterior cov / error correlation

# Thank you!

**Online supplementary material**

*github.com/lionel-rigoux/tutorial-bayesian-inference*

- interactive app
- code of all algorithms
- selected references

**VBA-Toolbox**

*mbb-team.github.io/VBA-toolbox*

**Easy and reproducible writing workflow**

*pandemics.gitlab.io*