

Bayesian Inference

Thomas Nichols

With thanks

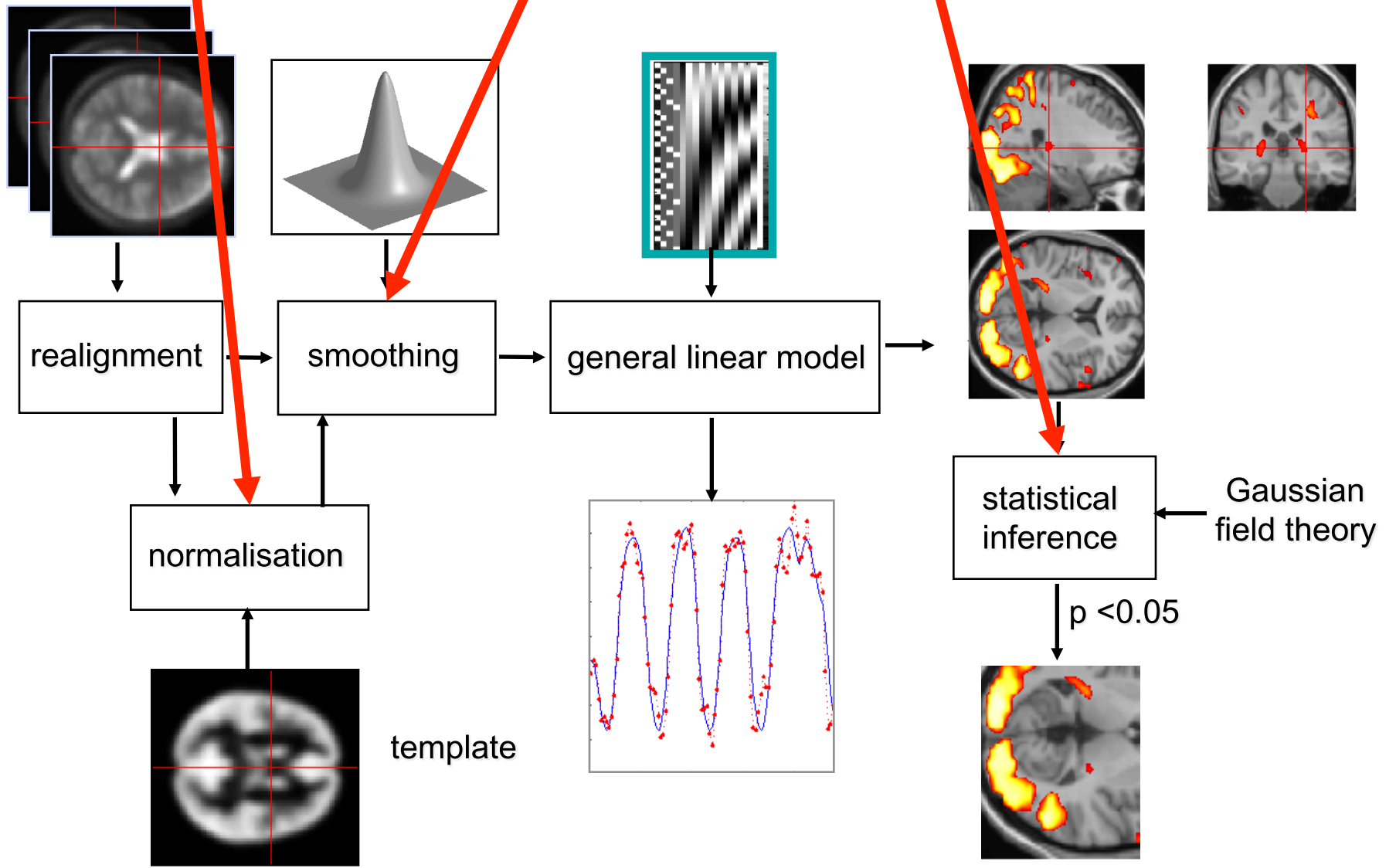
Lee Harrison

Bayesian segmentation and normalisation

Spatial priors on activation extent

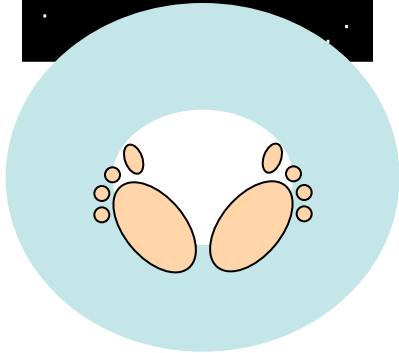
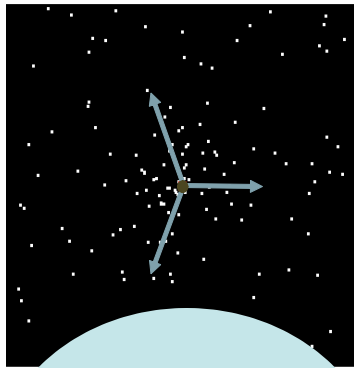
Posterior probability maps (PPMs)

Dynamic Causal Modelling



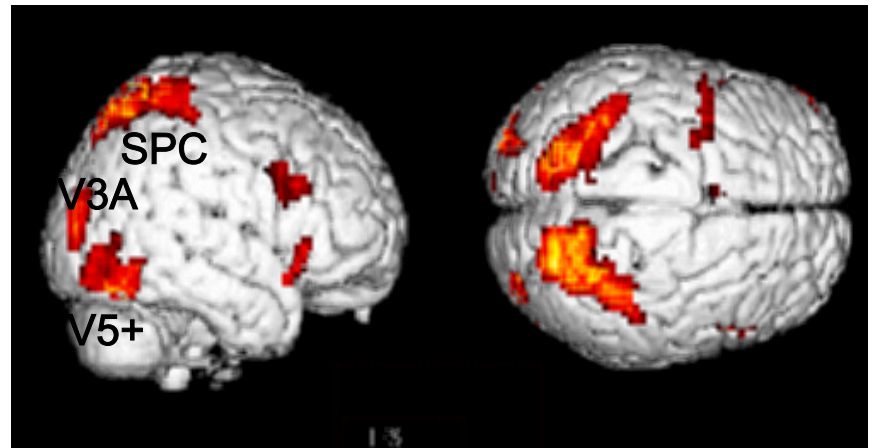
Attention to Motion

Paradigm



- fixation only
 - observe static dots
 - observe moving dots
 - task on moving dots
- + photic
 - + motion
 - + attention

Results



Attention – No attention

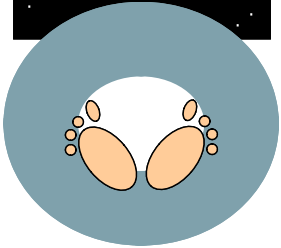
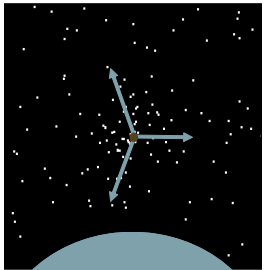
Büchel & Friston 1997, Cereb. Cortex

Büchel et al. 1998, Brain

- V1
- V5
- V5 + parietal cortex

Attention to Motion

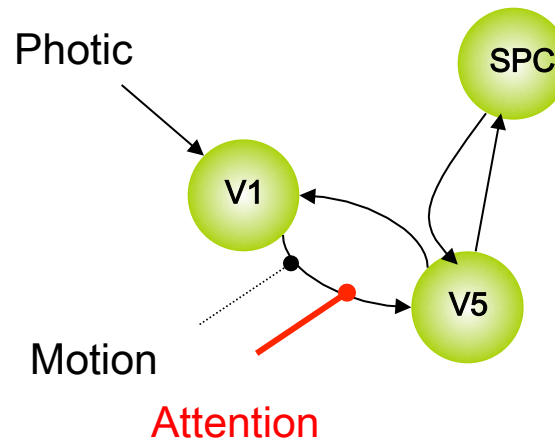
Paradigm



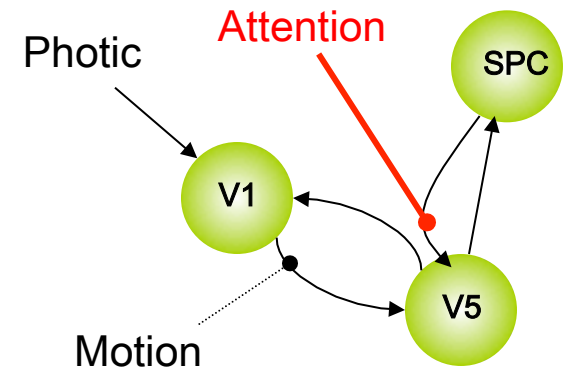
- fixation only
- observe static dots
- observe moving dots
- task on moving dots

Dynamic Causal Models

Model 1 (forward):
attentional modulation
of $V1 \rightarrow V5$: forward

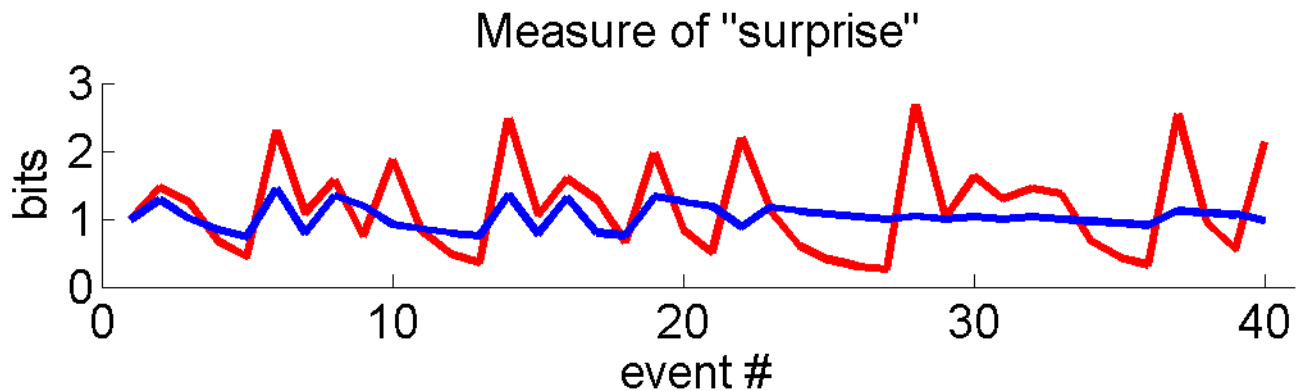
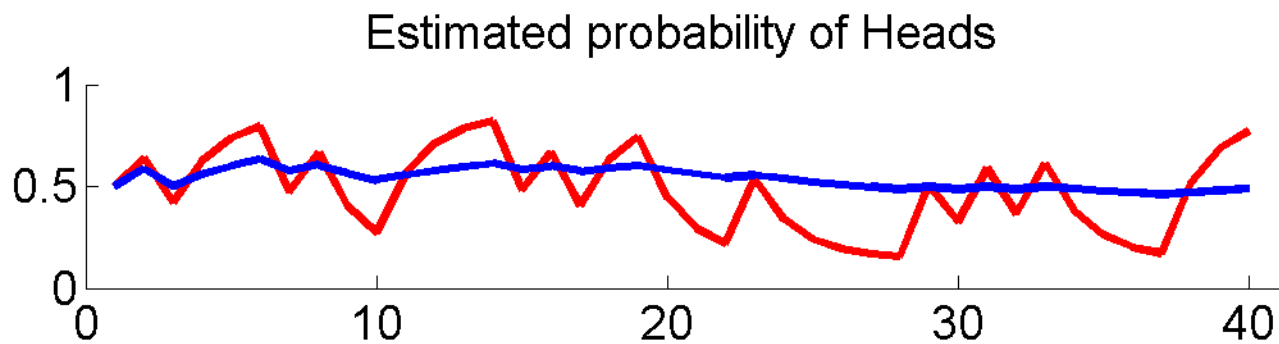
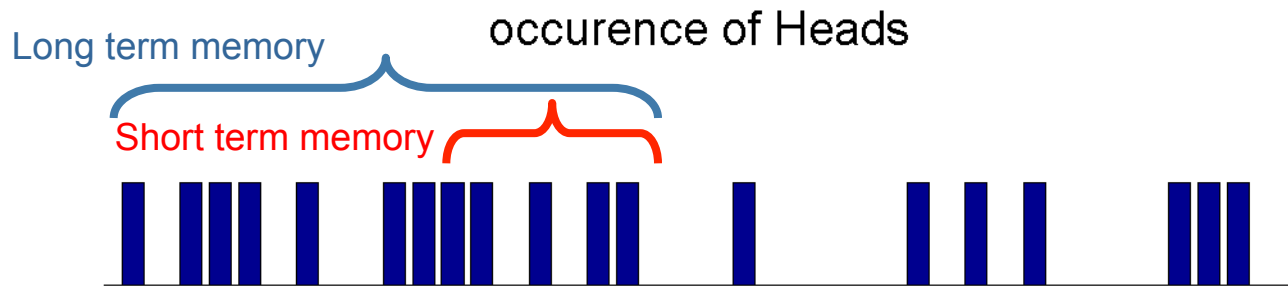


Model 2 (backward):
attentional modulation
of $SPC \rightarrow V5$: backward



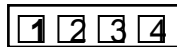
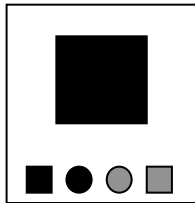
Bayesian model selection: Which model is optimal?

Responses to Uncertainty



Responses to Uncertainty

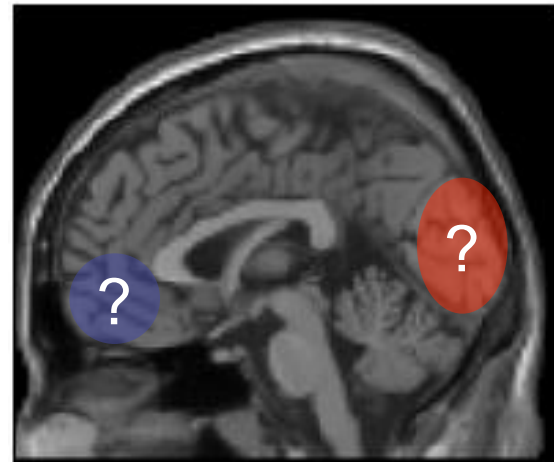
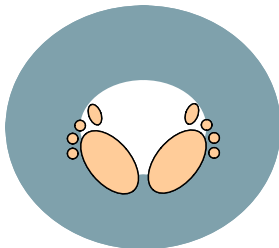
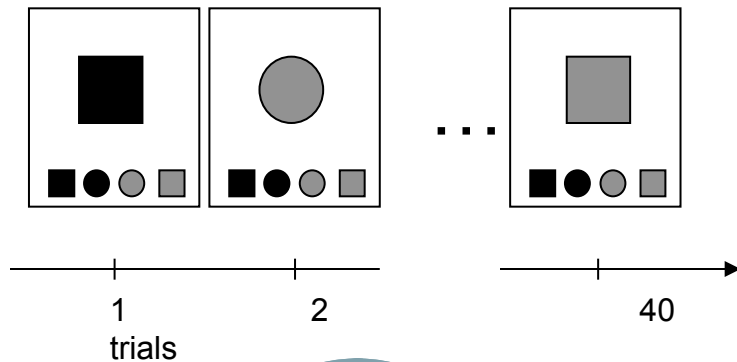
Paradigm



Stimuli sequence of randomly sampled discrete events

Model simple computational model of an observers response to uncertainty based on the number of past events (extent of memory)

Question which regions are best explained by short / long term memory model?



Overview

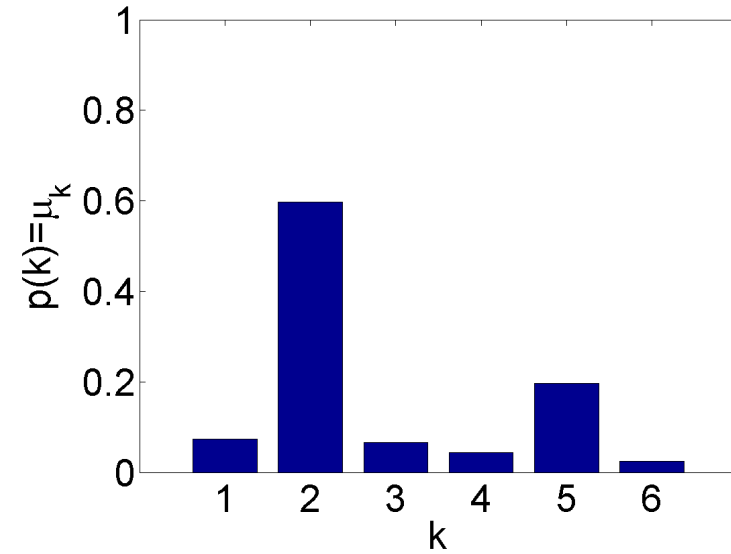
- Introductory remarks
- Some probability densities/distributions
- Probabilistic (generative) models
- Bayesian inference
- A simple example – Bayesian linear regression
- SPM applications
 - Segmentation
 - Dynamic causal modeling
 - Spatial models of fMRI time series

Probability distributions and densities

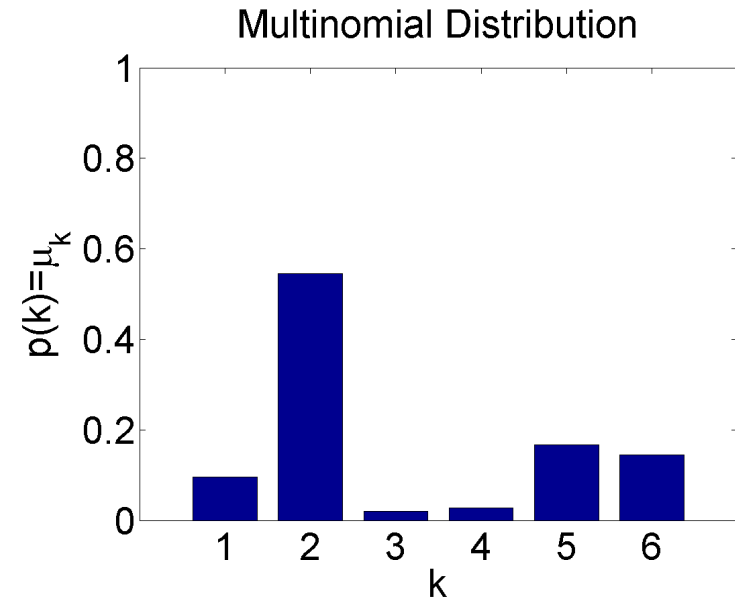


k=2

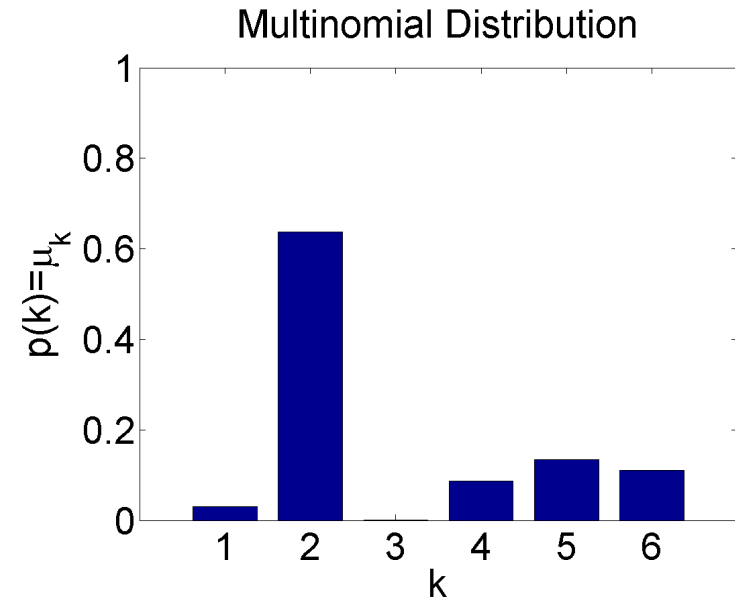
Multinomial Distribution



Probability distributions and densities



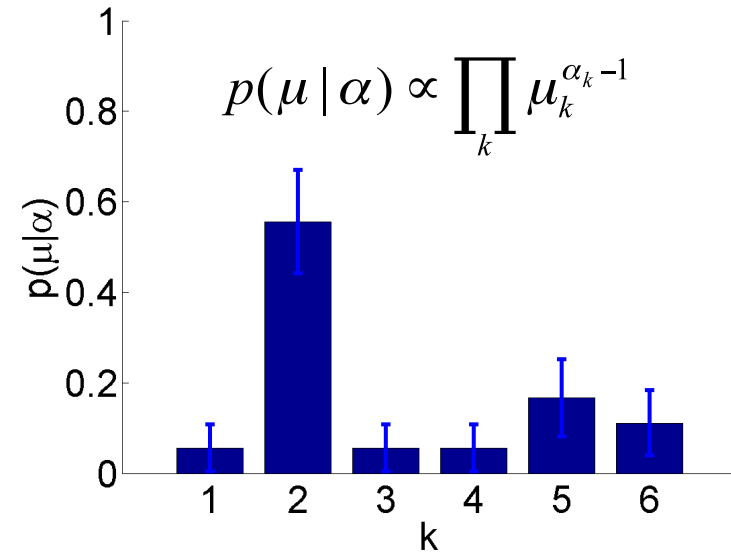
Probability distributions and densities



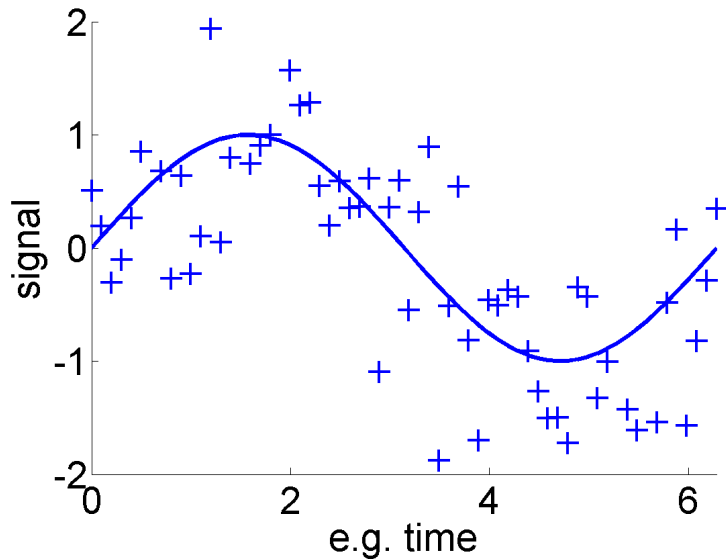
Probability distributions and densities



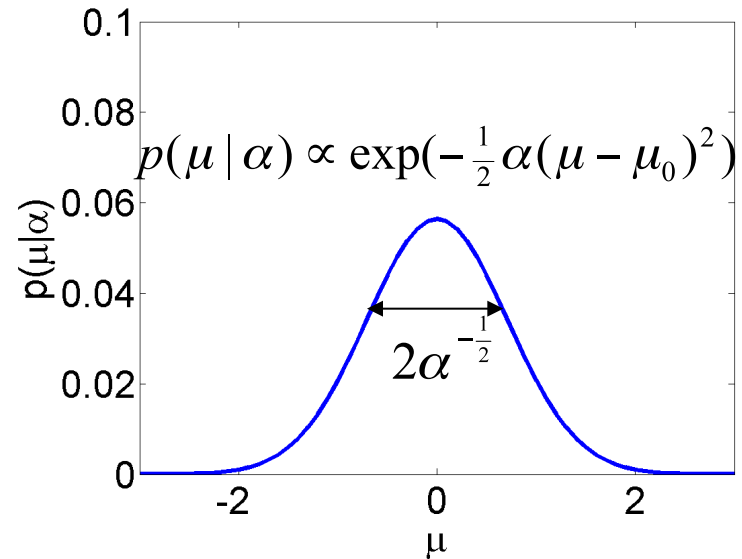
Dirichlet Distribution



Sine + Gaussian Noise



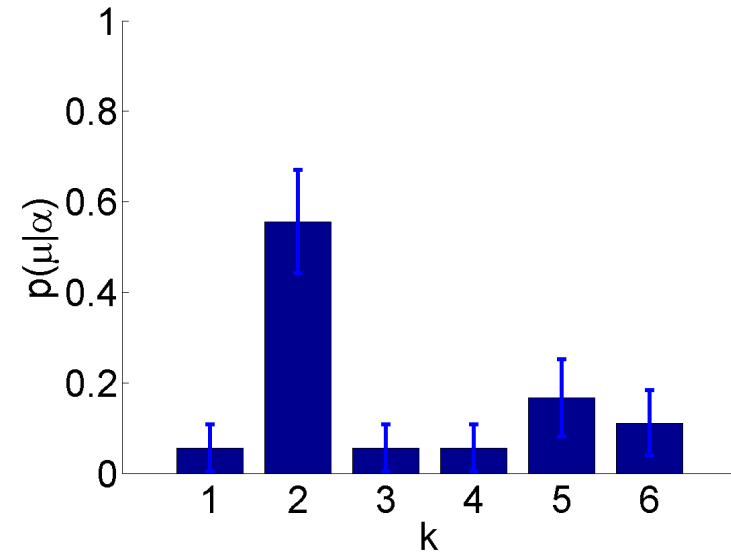
Gaussian Distribution



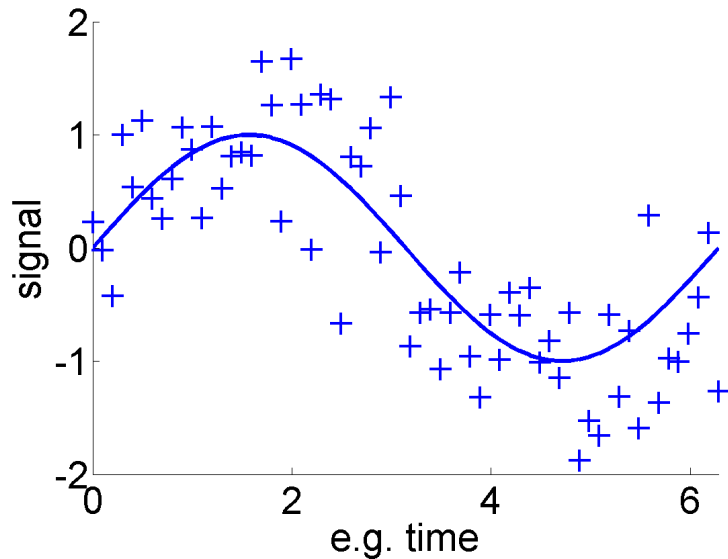
Probability distributions and densities



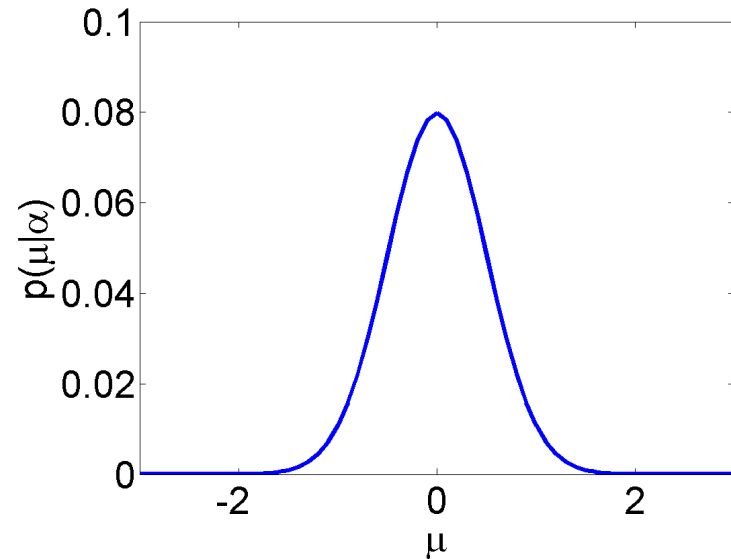
Dirichlet Distribution



Sine + Gaussian Noise



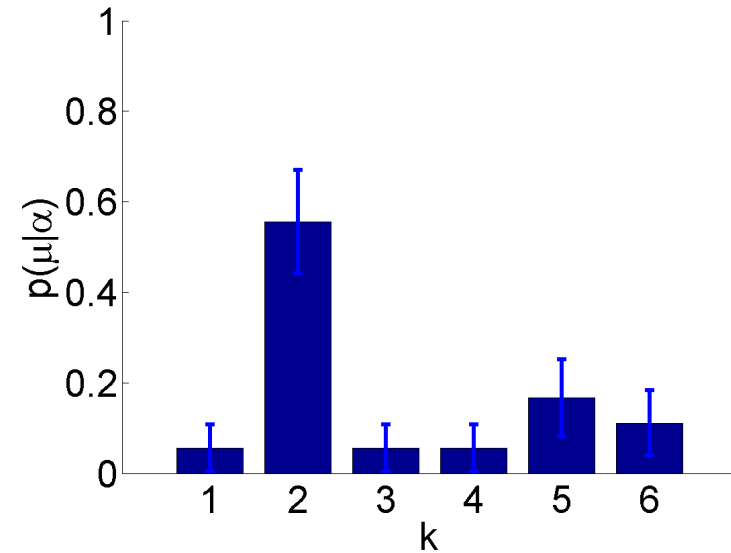
Gaussian Distribution



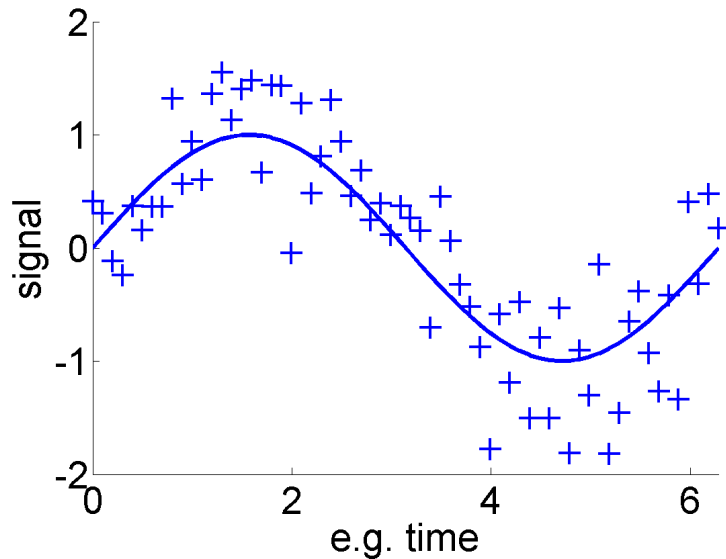
Probability distributions and densities



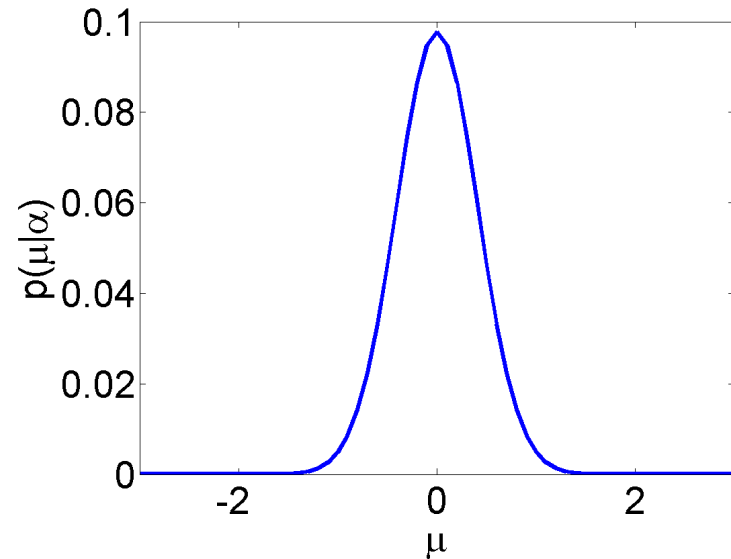
Dirichlet Distribution



Sine + Gaussian Noise



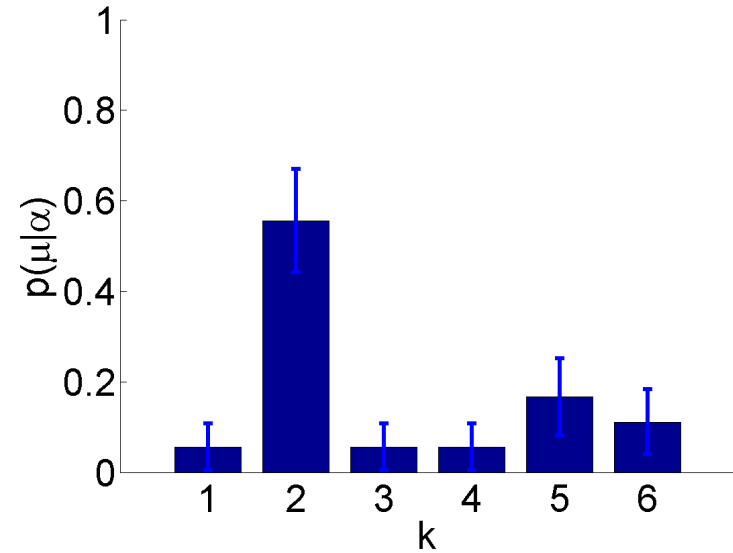
Gaussian Distribution



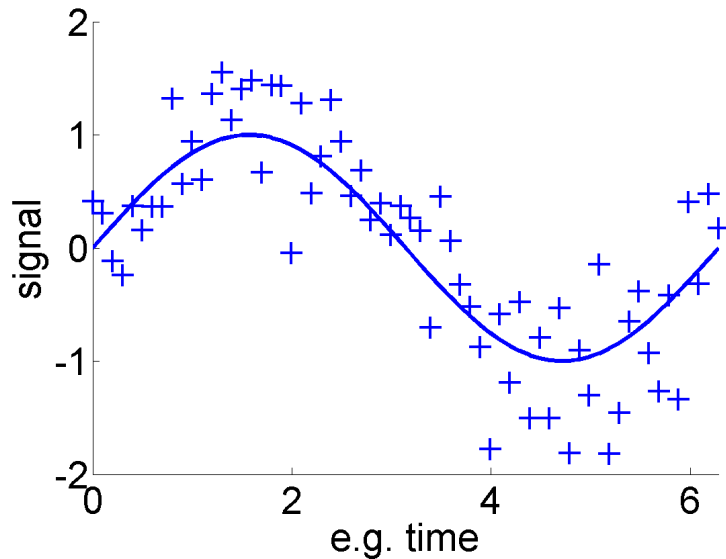
Probability distributions and densities



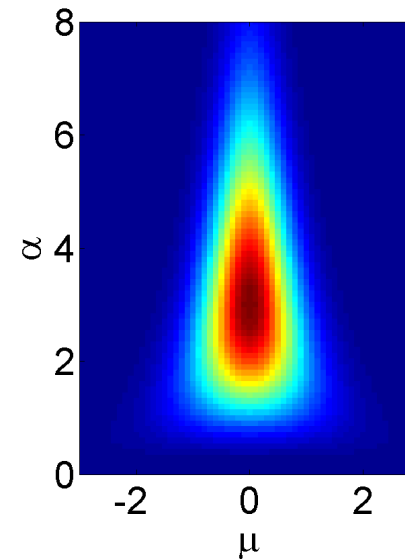
Dirichlet Distribution



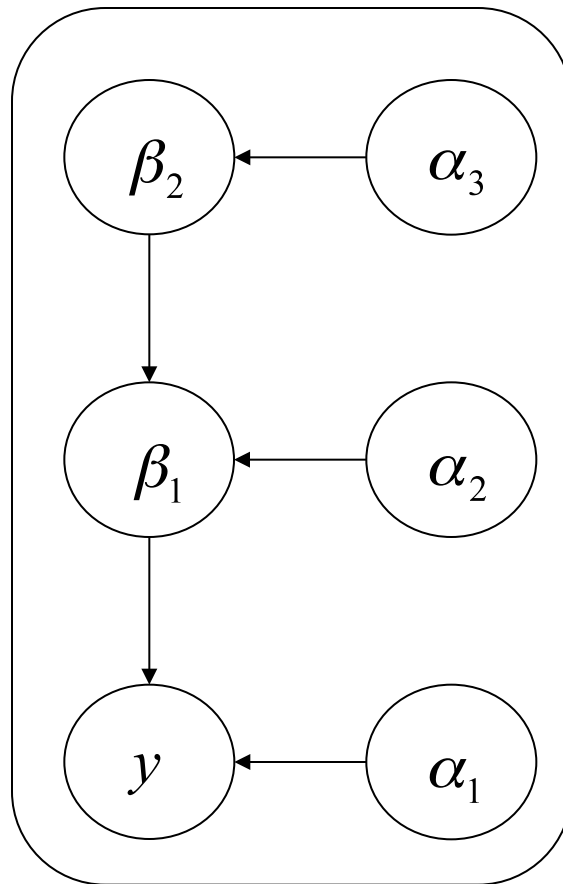
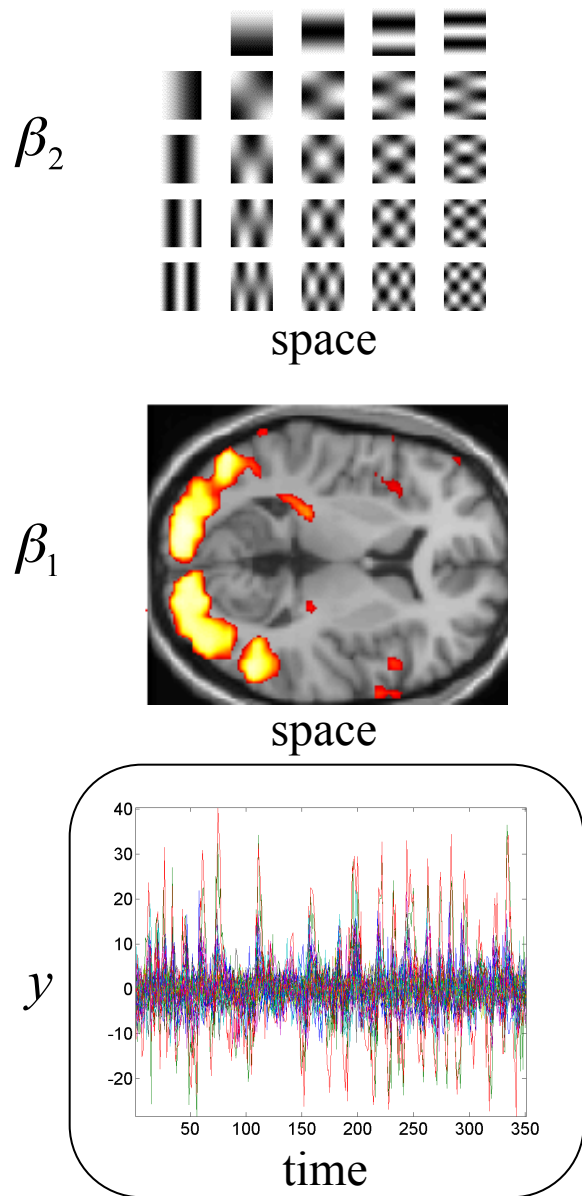
Sine + Gaussian Noise



Gaussian-Gamma Distribution



Generative models



$$\theta = \{\beta, \alpha\}$$

$$\beta_2 = e_3 \sim N(0, \alpha_3^{-1})$$

$$\beta_1 = X_2 \beta_2 + e_2$$

$$y = X_1 \beta_1 + e_1$$

generation

$q(\theta)$?
estimation



Bayesian statistics

new data

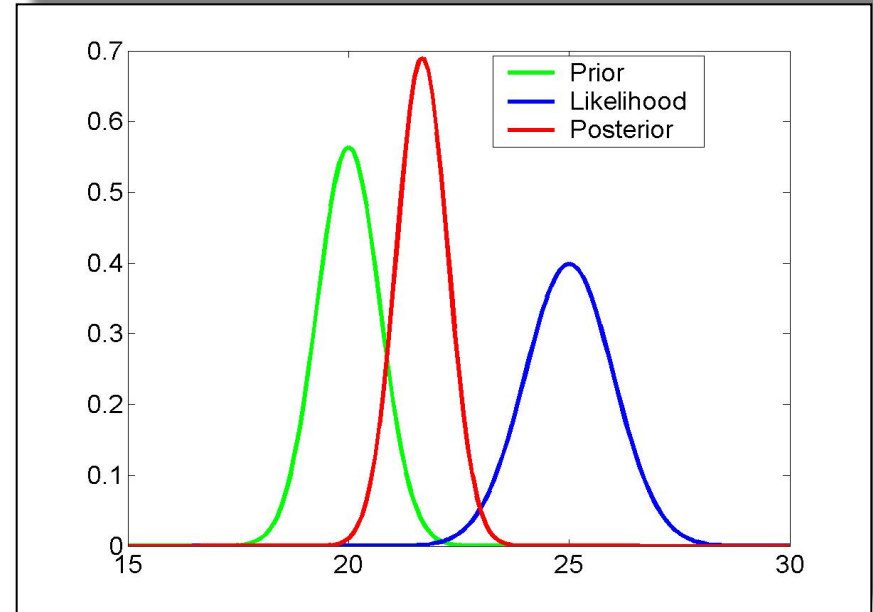
$$p(y | \theta)$$

prior knowledge

$$p(\theta)$$

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

posterior \propto likelihood \cdot prior



Bayes theorem allows one to formally incorporate prior knowledge into computing statistical probabilities.

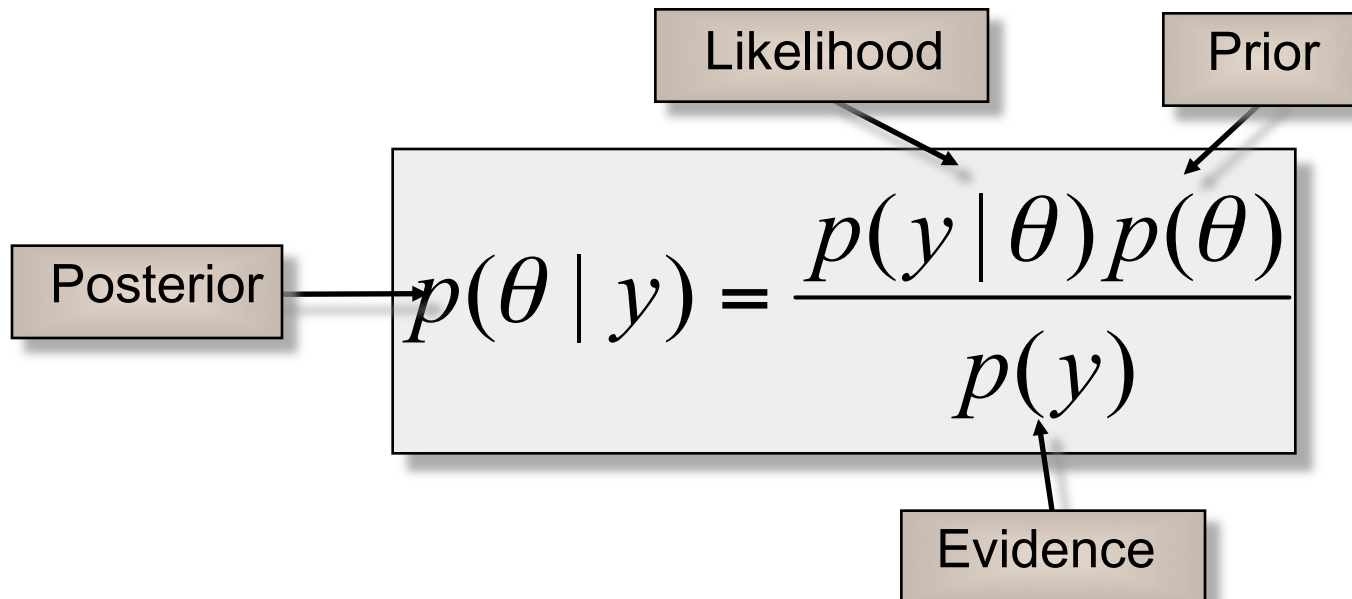
The “posterior” probability of the parameters given the data is an optimal combination of prior knowledge and new data, weighted by their relative precision.

Bayes' rule

Given data y and parameters θ , their joint probability can be written in 2 ways:

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

Eliminating $p(y, \theta)$ gives Bayes' rule:



Principles of Bayesian inference

⇒ Formulation of a generative model

likelihood $p(y|\theta)$
prior distribution $p(\theta)$

⇒ Observation of data

y

⇒ Update of beliefs based upon observations, given a prior state of knowledge

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

Univariate Gaussian

Normal densities

$$p(\beta) = N(\beta; \mu_p, \alpha_p^{-1})$$

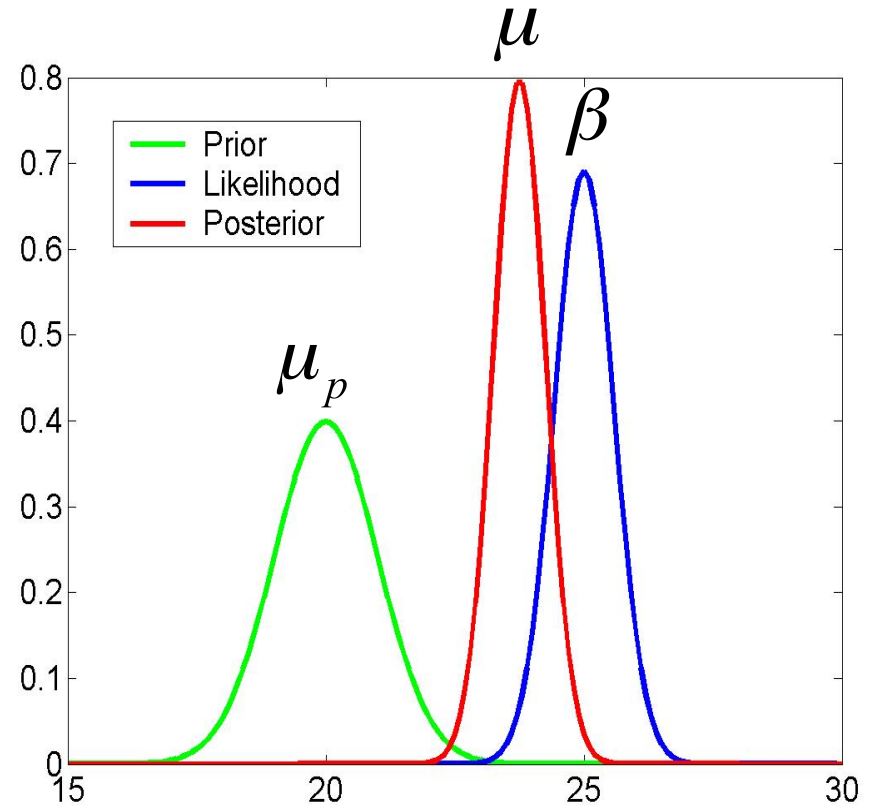
$$p(y | \beta) = N(y; \beta, \alpha_e^{-1})$$

$$p(\beta | y) = N(\beta; \mu, \alpha^{-1})$$

$$\alpha = \alpha_e + \alpha_p$$
$$\mu = \alpha^{-1} (\alpha_e y + \alpha_p \mu_p)$$

$$y = \beta + e$$

Posterior mean =
precision-weighted combination of
prior mean and data mean



Bayesian GLM: univariate case

Normal densities

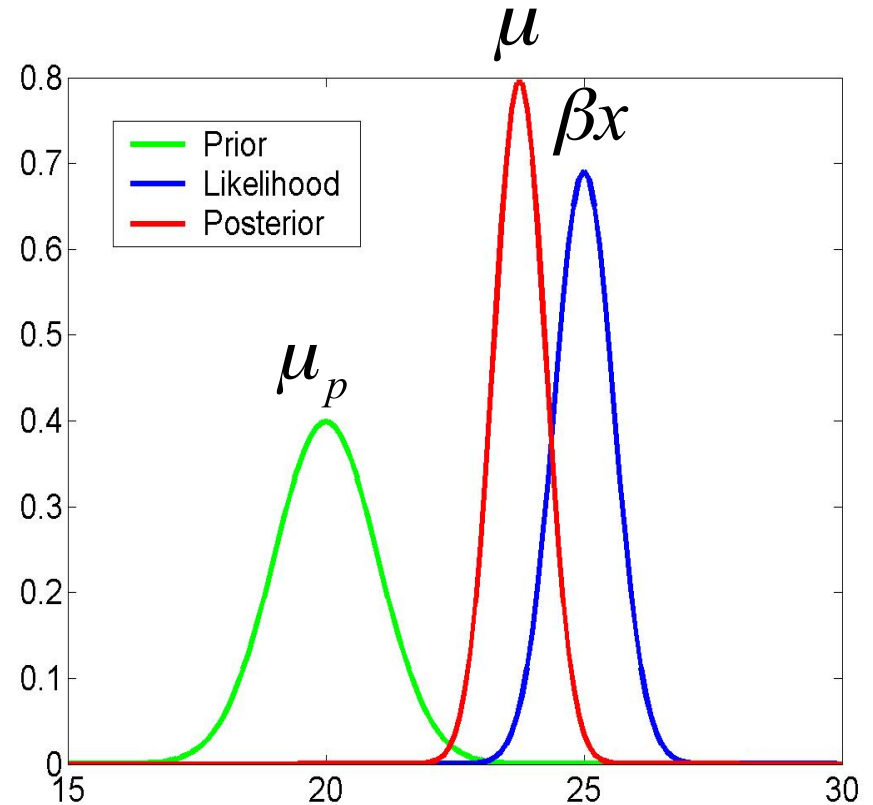
$$p(\beta) = N(\beta; \mu_p, \alpha_p^{-1})$$

$$p(y | \beta) = N(y; \beta x, \alpha_e^{-1})$$

$$p(\beta | y) = N(\beta; \mu, \alpha^{-1})$$

$$\alpha = \alpha_e x^2 + \alpha_p$$
$$\mu = \alpha^{-1} (\alpha_e x y + \alpha_p \mu_p)$$

$$y = \beta x + e$$



Bayesian GLM: multivariate case

Normal densities

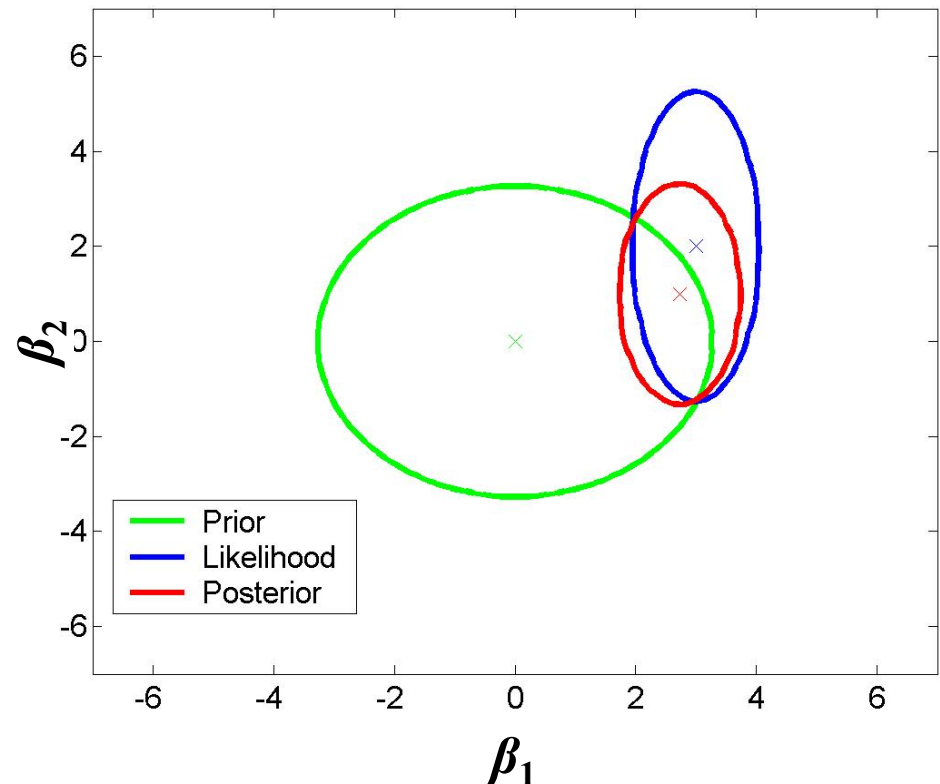
$$p(\beta) = N(\beta; \mu_p, C_p)$$

$$p(y | \beta) = N(y; X\beta, C_e)$$

$$p(\beta | y) = N(\beta; \mu, C)$$

$$C^{-1} = X^T C_e^{-1} X + C_p^{-1}$$
$$\mu = C(X^T C_e y + C_p^{-1} \mu_p)$$

$$y = X\beta + e$$

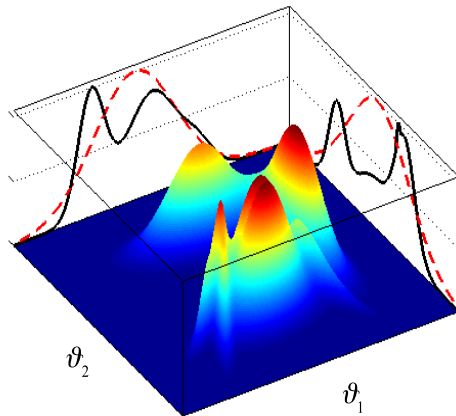


One step if C_e and C_p are known.
Otherwise iterative estimation.

Approximate inference: optimization

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

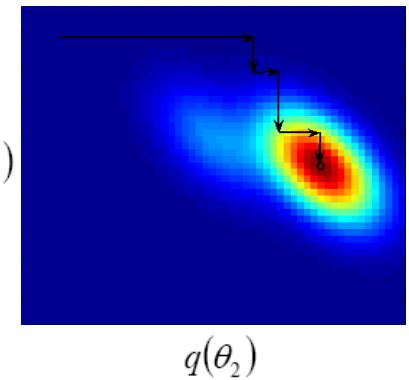
mean-field approximation



Approximate posterior

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

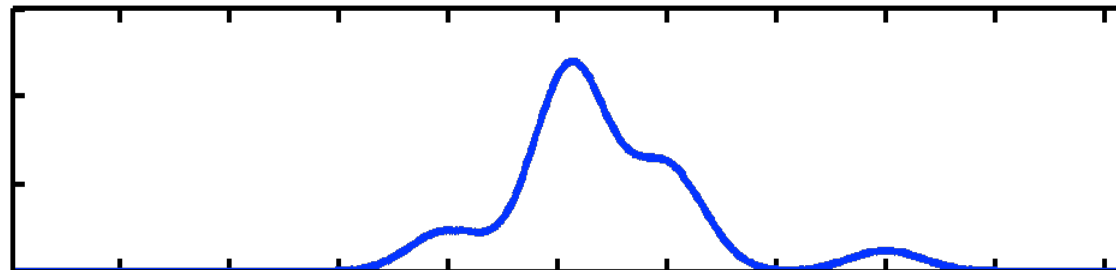
iteratively improve



free energy

$$\log p(y|m) = \int_{\theta} q(\theta) \log \left(\frac{p(y, \theta | m)}{q(\theta)} \right) + \int_{\theta} q(\theta) \log \left(\frac{q(\theta)}{p(\theta | y, m)} \right)$$

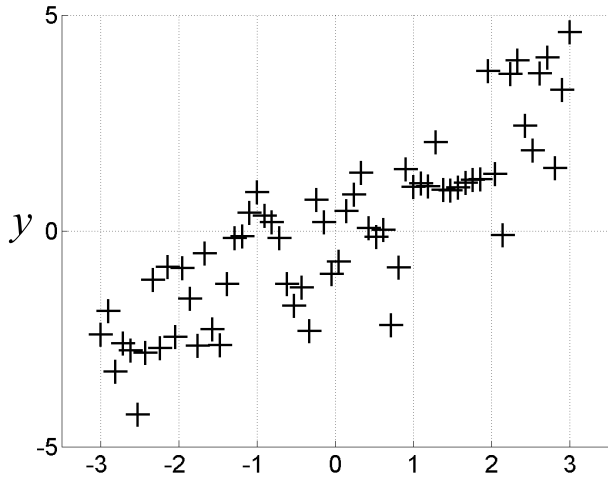
Objective function



Value of parameter

Simple example – linear regression

Data



Ordinary least squares

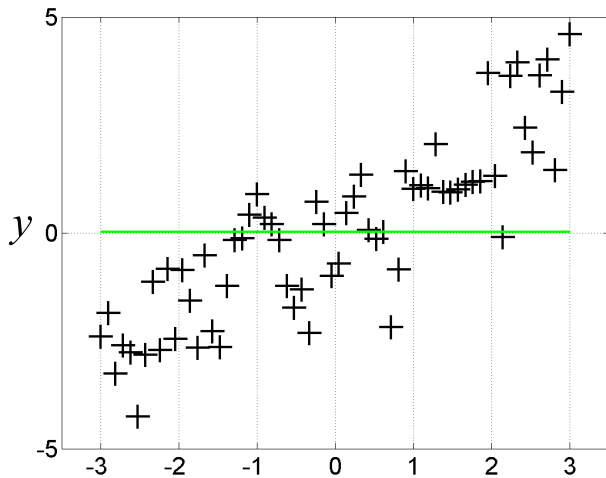
$$y = X\beta$$

$$E_D = (y - X\beta)^T (y - X\beta)$$

$$\frac{\partial E_D}{\partial \beta} = 0 \Rightarrow \hat{\beta}_{ols} = (X^T X)^{-1} X^T y$$

Simple example – linear regression

Data and model fit



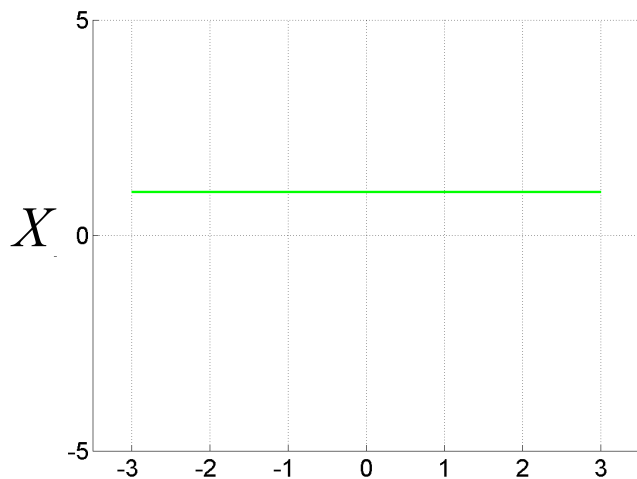
Ordinary least squares

$$y = X\beta$$

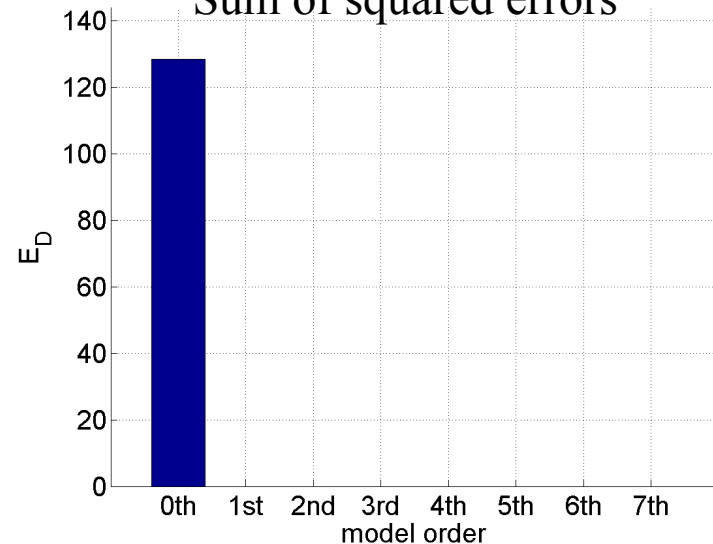
$$E_D = (y - X\beta)^T (y - X\beta)$$

$$\frac{\partial E_D}{\partial \beta} = 0 \Rightarrow \hat{\beta}_{ols} = (X^T X)^{-1} X^T y$$

Bases (explanatory variables)

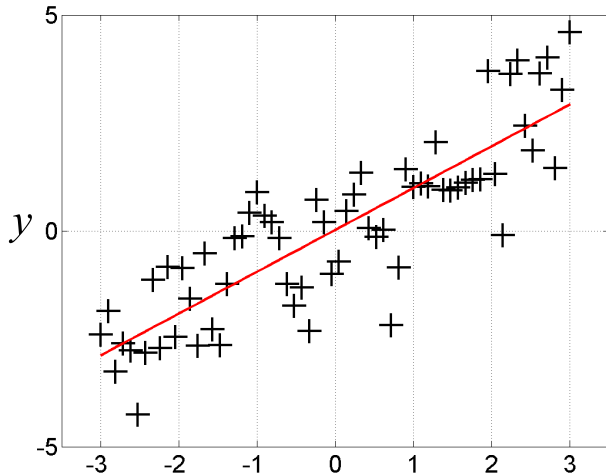


Sum of squared errors



Simple example – linear regression

Data and model fit



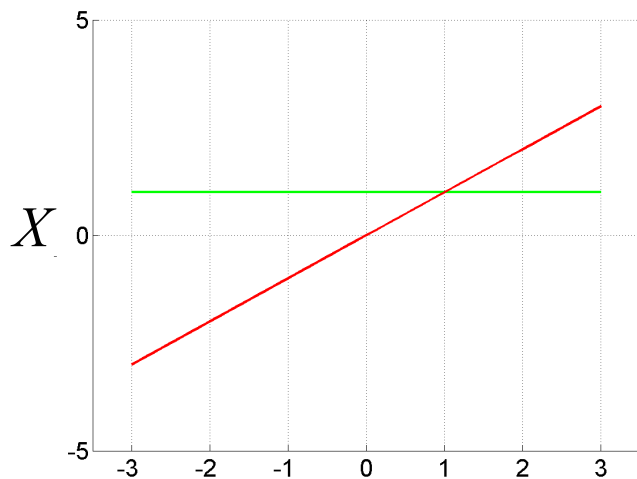
Ordinary least squares

$$y = X\beta$$

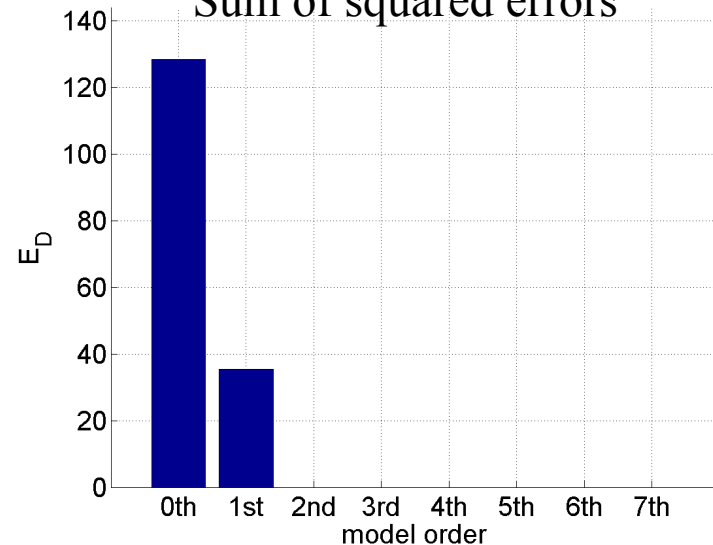
$$E_D = (y - X\beta)^T (y - X\beta)$$

$$\frac{\partial E_D}{\partial \beta} = 0 \Rightarrow \hat{\beta}_{ols} = (X^T X)^{-1} X^T y$$

Bases (explanatory variables)

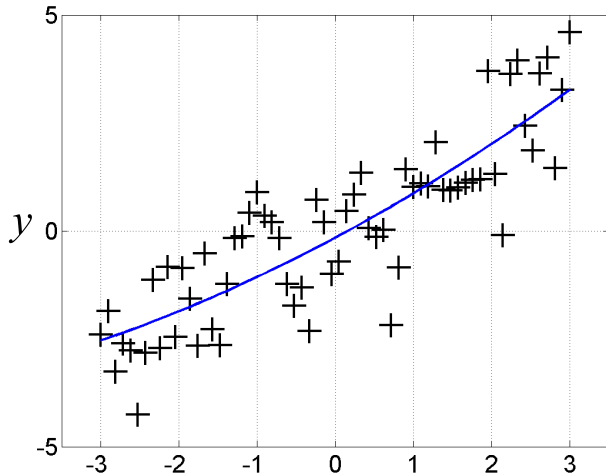


Sum of squared errors



Simple example – linear regression

Data and model fit



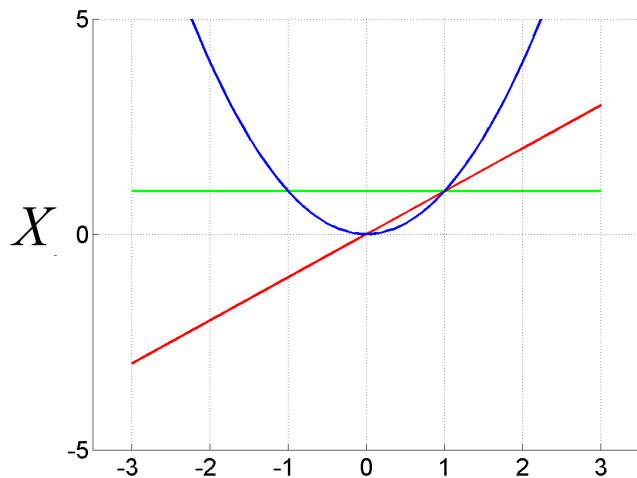
Ordinary least squares

$$y = X\beta$$

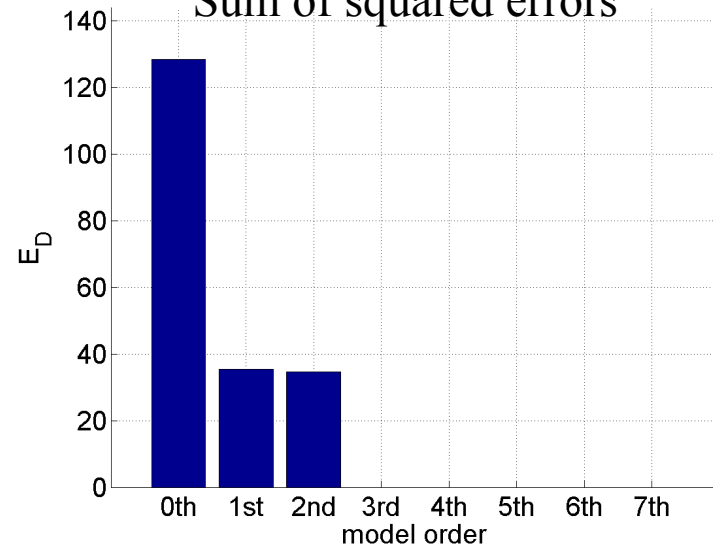
$$E_D = (y - X\beta)^T (y - X\beta)$$

$$\frac{\partial E_D}{\partial \beta} = 0 \Rightarrow \hat{\beta}_{ols} = (X^T X)^{-1} X^T y$$

Bases (explanatory variables)

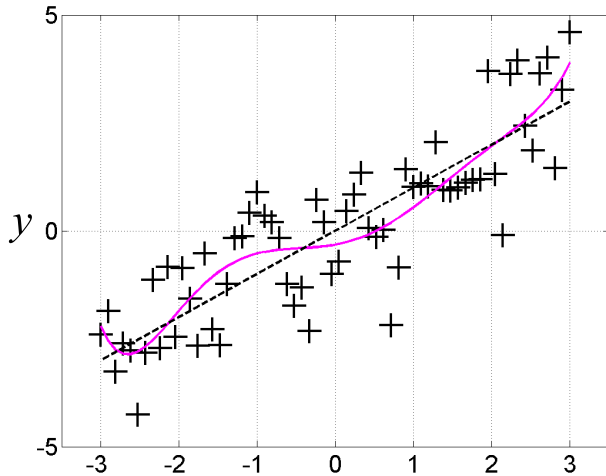


Sum of squared errors



Simple example – linear regression

Data and model fit



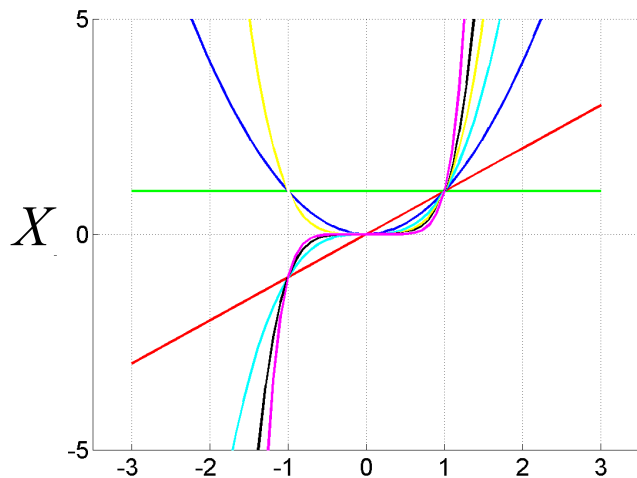
Ordinary least squares

Over-fitting: model fits noise

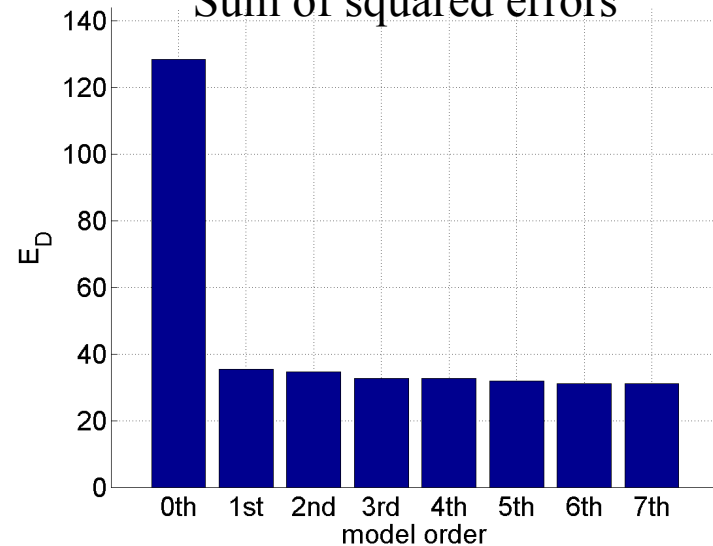
Inadequate cost function: blind to overly complex models

Solution: include uncertainty in model parameters

Bases (explanatory variables)

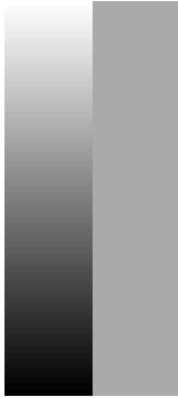


Sum of squared errors



Bayesian linear regression: *priors and likelihood*

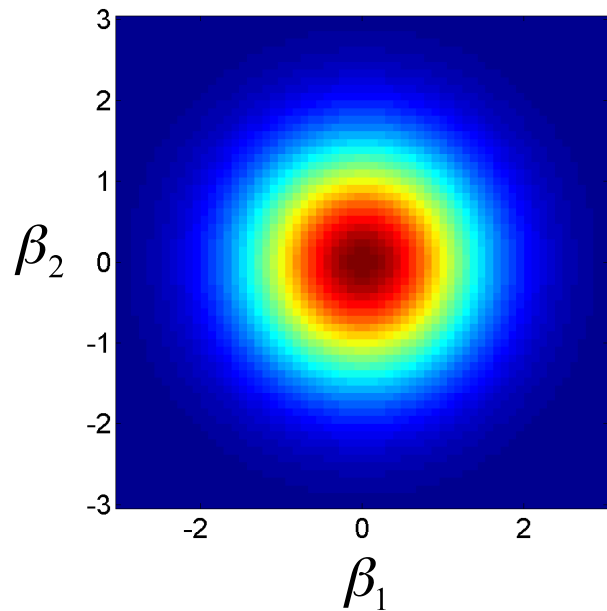
$X =$



Model:

$$y = X\beta + e$$

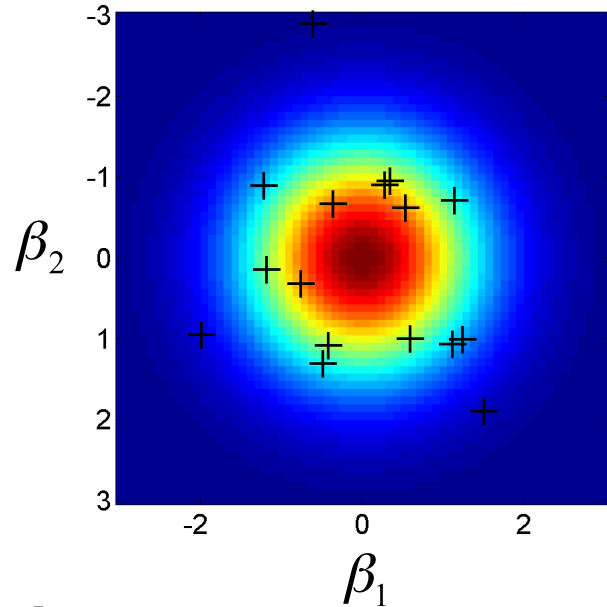
Bayesian linear regression: *priors and likelihood*



Model: $y = X\beta + e$

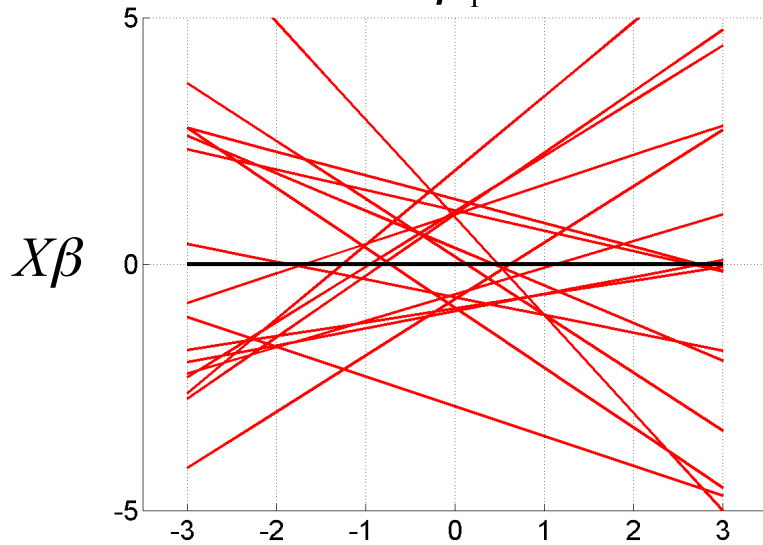
Prior: $p(\beta|\alpha_2) = N_k(0, \alpha_2^{-1}I_k)$
 $\propto \exp(-\alpha_2\|\beta\|^2 / 2)$



Bayesian linear regression: *priors and likelihood*



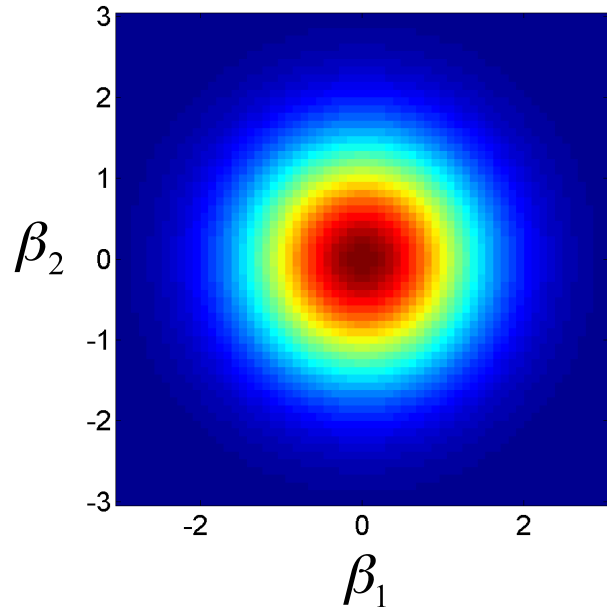
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$



-  Sample curves from prior
(before observing any data)
-  Mean curve

Bayesian linear regression: *priors and likelihood*



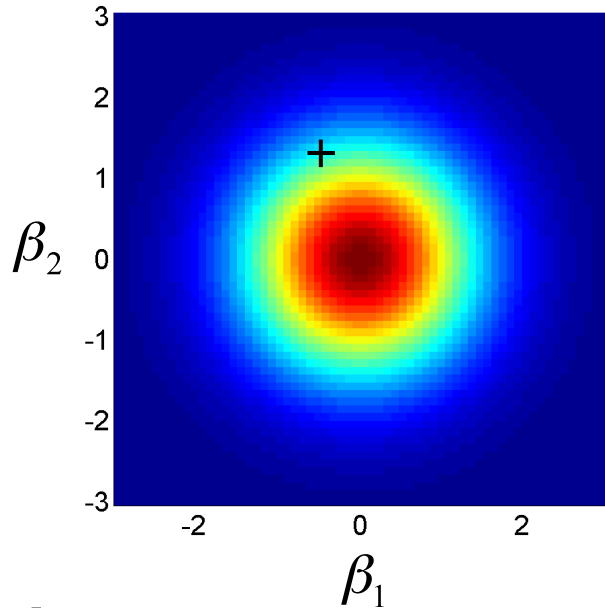
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

Likelihood:

$$p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1^{-1})$$
$$p(y_i | \beta, \alpha_1) = N(X_i \beta, \alpha_1^{-1})$$
$$\propto \exp(-\alpha_1 (y_i - X_i \beta)^2 / 2)$$

Bayesian linear regression: *priors and likelihood*



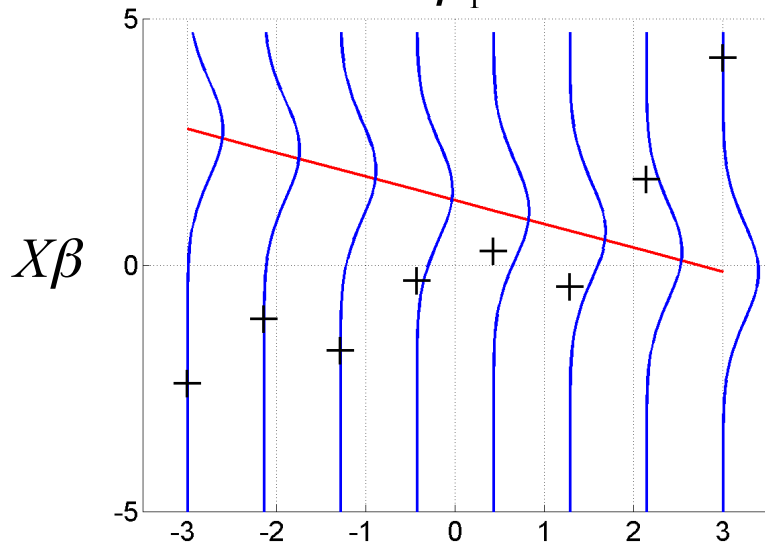
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

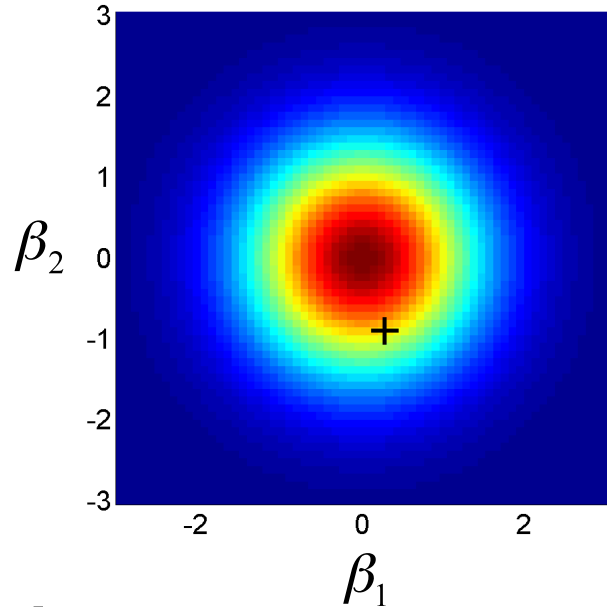
Likelihood:

$$p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1^{-1})$$

$$p(y_i | \beta, \alpha_1) = N(X_i \beta, \alpha_1^{-1})$$
$$\propto \exp(-\alpha_1 (y_i - X_i \beta)^2 / 2)$$



Bayesian linear regression: *priors and likelihood*



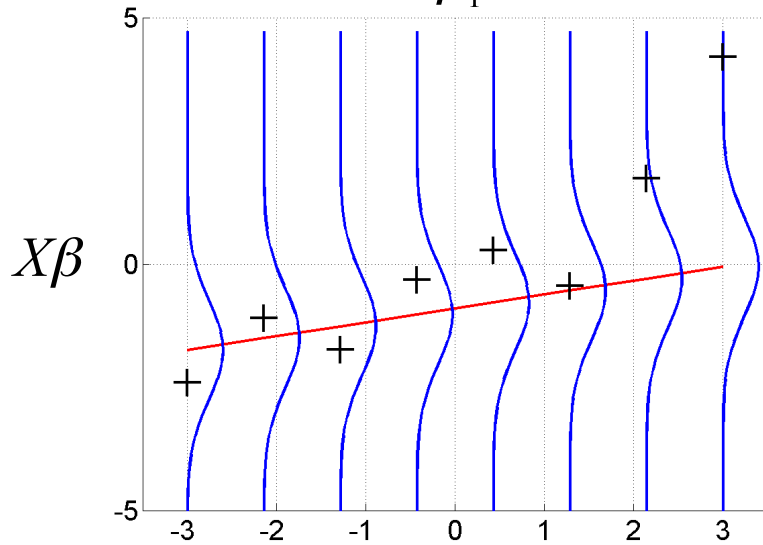
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

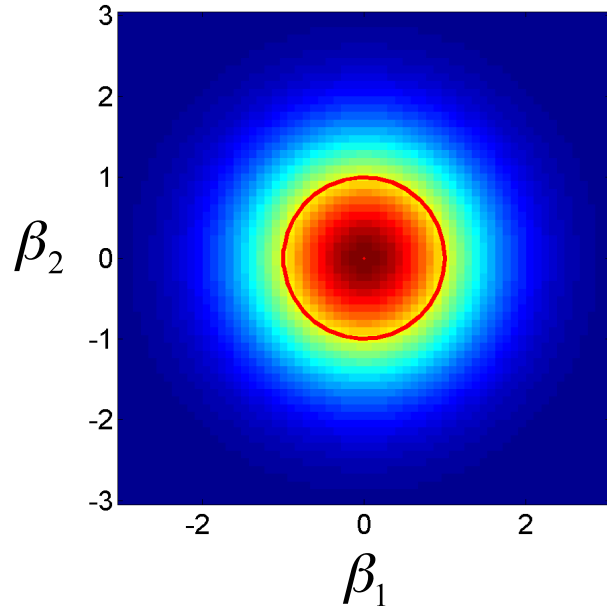
Likelihood:

$$p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1^{-1})$$

$$p(y_i | \beta, \alpha_1) = N(X_i \beta, \alpha_1^{-1})$$
$$\propto \exp(-\alpha_1 (y_i - X_i \beta)^2 / 2)$$



Bayesian linear regression: *posterior*

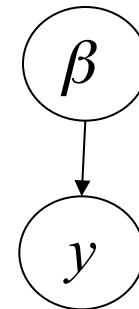


Model: $y = X\beta + e$

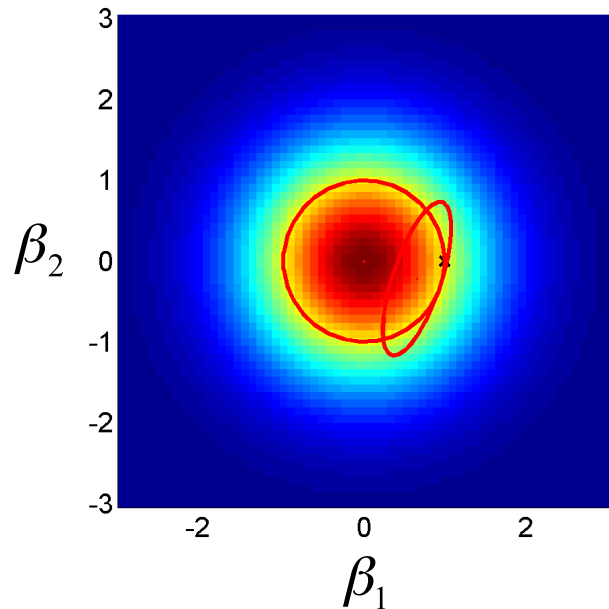
Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

Likelihood: $p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1)$

Bayes Rule: $p(\beta | y, \alpha) \propto p(y | \beta, \alpha) p(\beta | \alpha)$



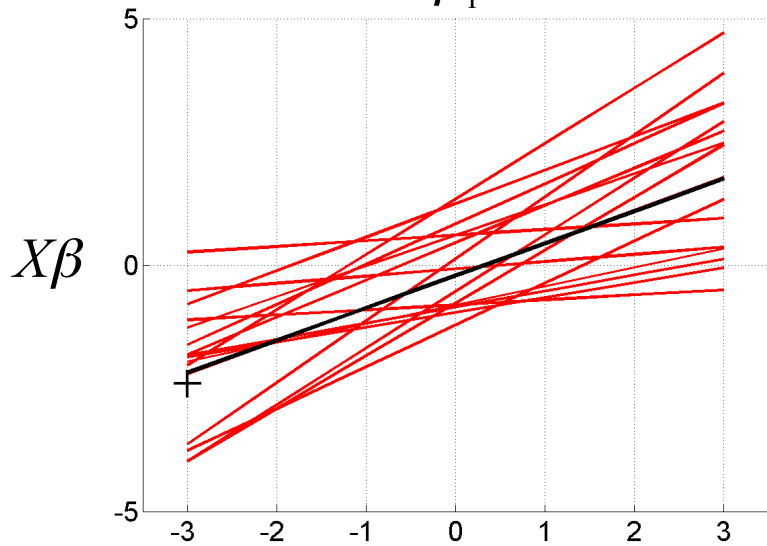
Bayesian linear regression: *posterior*



Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

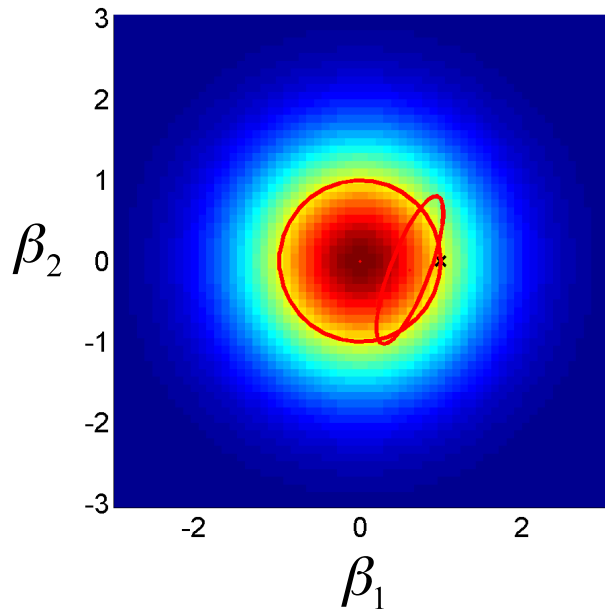
Likelihood: $p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1)$



Bayes Rule: $p(\beta | y, \alpha) \propto p(y | \beta, \alpha) p(\beta | \alpha)$

Posterior: $p(\beta | y, \alpha) = N(\mu, C)$
 $C = (\alpha_1 X^T X + \alpha_2 I_k)^{-1}$
 $\mu = \alpha_1 C X^T y$

Bayesian linear regression: *posterior*



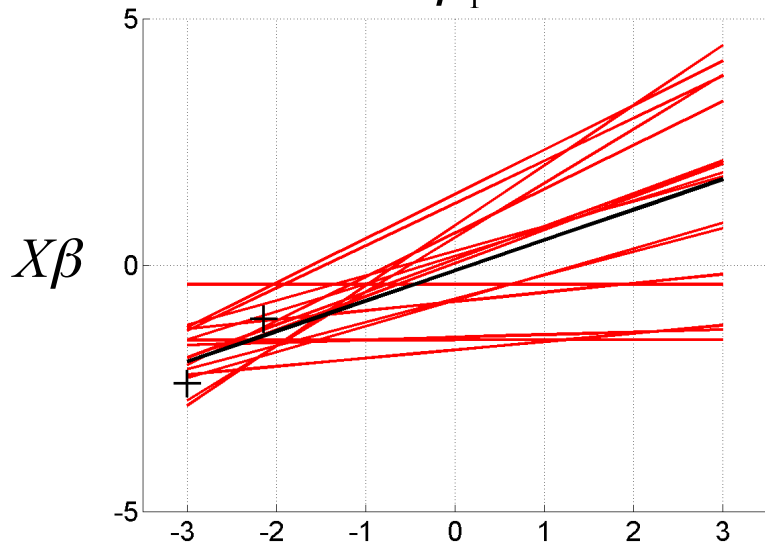
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

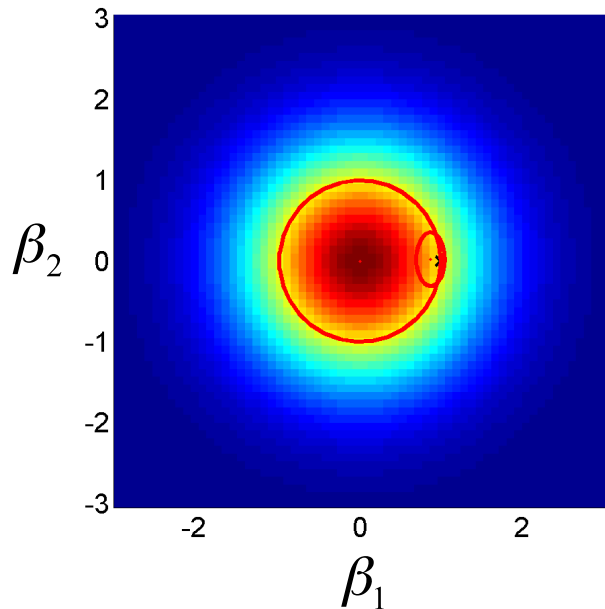
Likelihood: $p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1)$

Bayes Rule: $p(\beta | y, \alpha) \propto p(y | \beta, \alpha) p(\beta | \alpha)$

Posterior: $p(\beta | y, \alpha) = N(\mu, C)$
 $C = (\alpha_1 X^T X + \alpha_2 I_k)^{-1}$
 $\mu = \alpha_1 C X^T y$



Bayesian linear regression: *posterior*



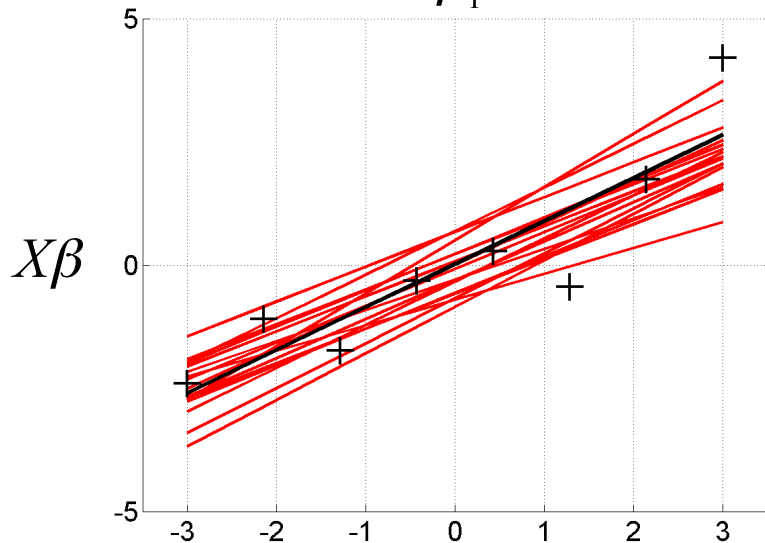
Model: $y = X\beta + e$

Prior: $p(\beta | \alpha_2) = N_k(0, \alpha_2^{-1} I_k)$
 $\propto \exp(-\alpha_2 \|\beta\|^2 / 2)$

Likelihood: $p(y | \beta, \alpha_1) = \prod_{i=1}^N p(y_i | \beta, \alpha_1)$

Bayes Rule: $p(\beta | y, \alpha) \propto p(y | \beta, \alpha) p(\beta | \alpha)$

Posterior: $p(\beta | y, \alpha) = N(\mu, C)$
 $C = (\alpha_1 X^T X + \alpha_2 I_k)^{-1}$
 $\mu = \alpha_1 C X^T y$



Posterior Probability Maps (PPMs)

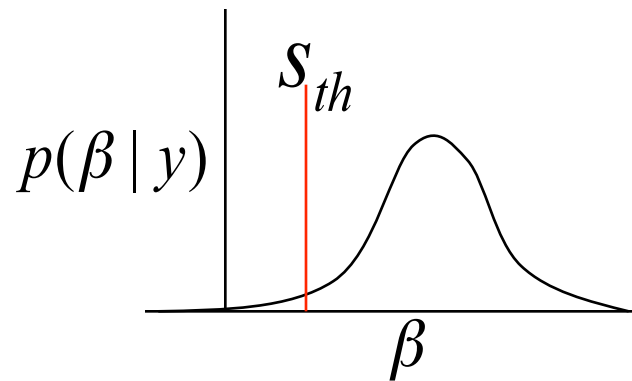
Posterior distribution: probability of the effect given the data

$$p(\beta | y)$$

mean: size of effect
precision: variability

Posterior probability map: images of the probability (confidence) that an activation exceeds some specified threshold s_{th} , given the data y

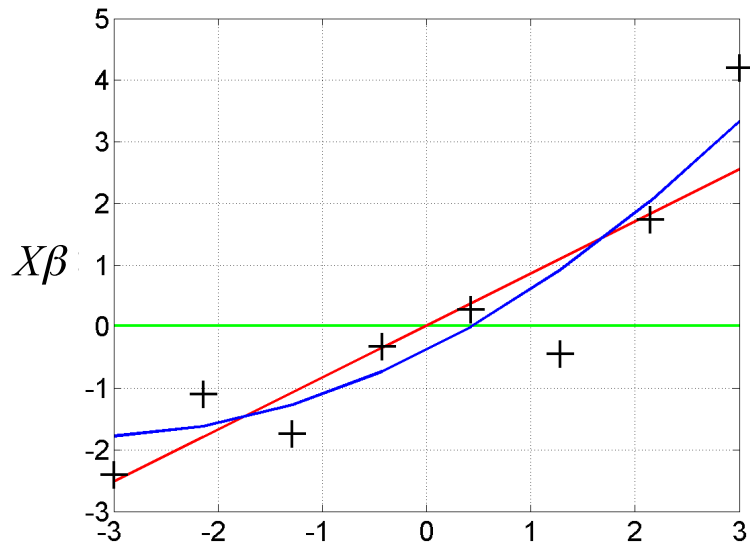
$$p(\beta > s_{th} | y) > p_{th}$$



Two thresholds:

- activation threshold s_{th} : percentage of whole brain mean signal (physiologically relevant size of effect)
- probability p_{th} that voxels must exceed to be displayed (e.g. 95%)

Bayesian linear regression: *model selection*



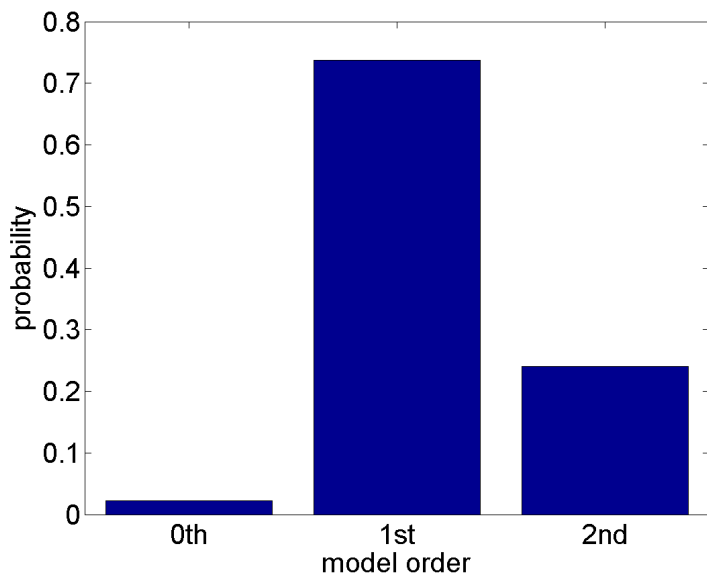
Bayes Rule:

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

normalizing constant

Model evidence:

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



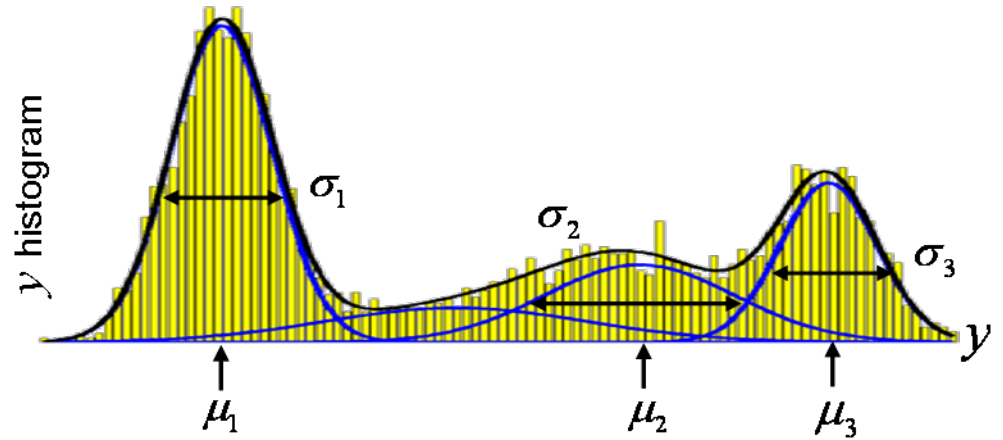
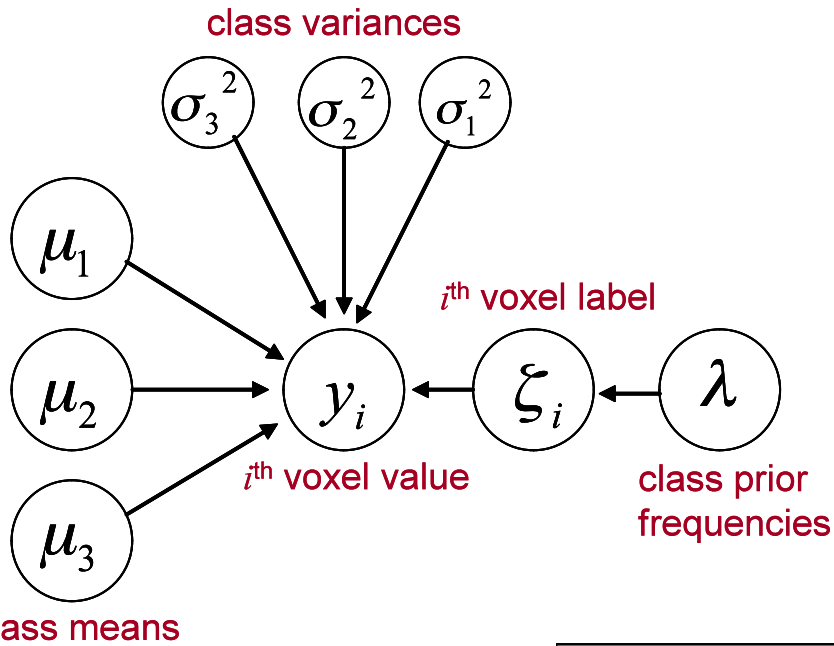
$$\log p(y | \alpha, m) =$$

$$accuracy(m) - complexity(m)$$

$$accuracy(m) \propto \|y - X\mu\|^2$$

$$complexity(m) \propto k \log \alpha_2^{-1} + \alpha_2 \|\mu\|^2$$

aMRI segmentation



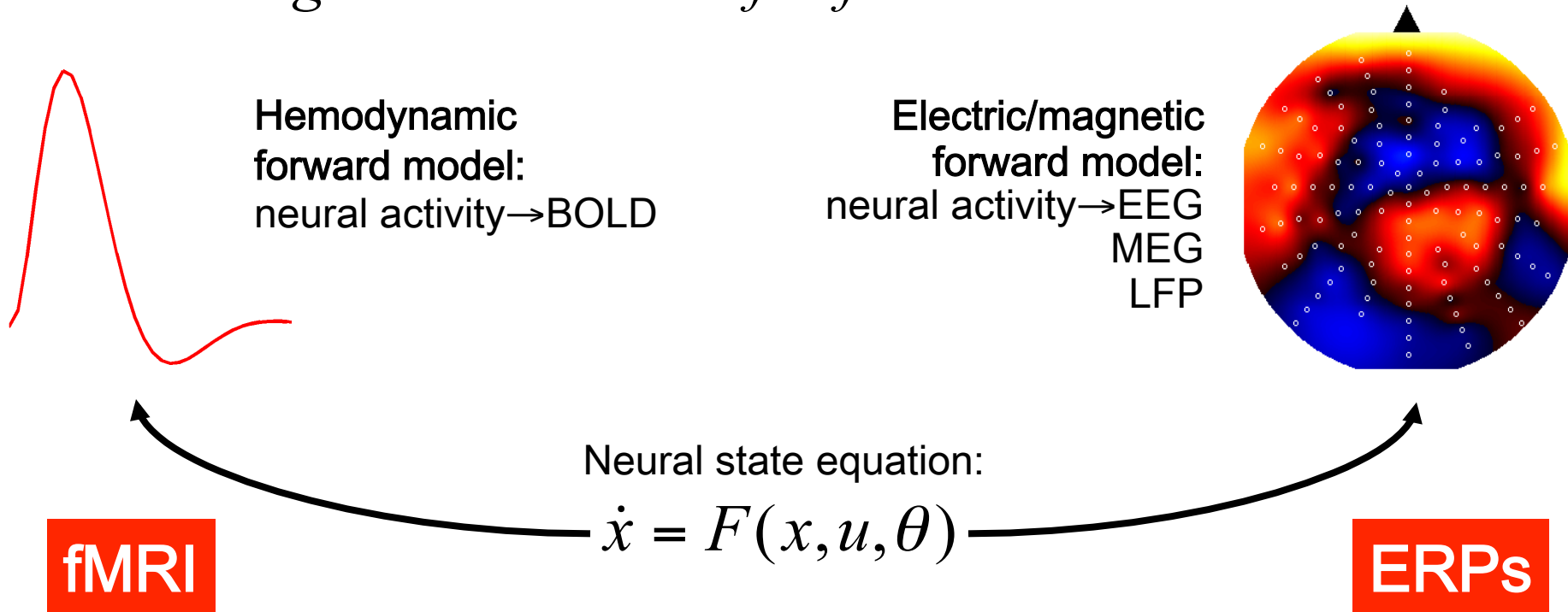
PPM of belonging to...

grey matter

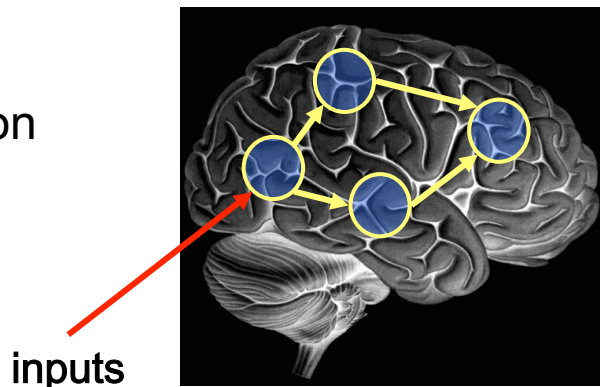
white matter

CSF

Dynamic Causal Modelling: *generative model for fMRI and ERPs*

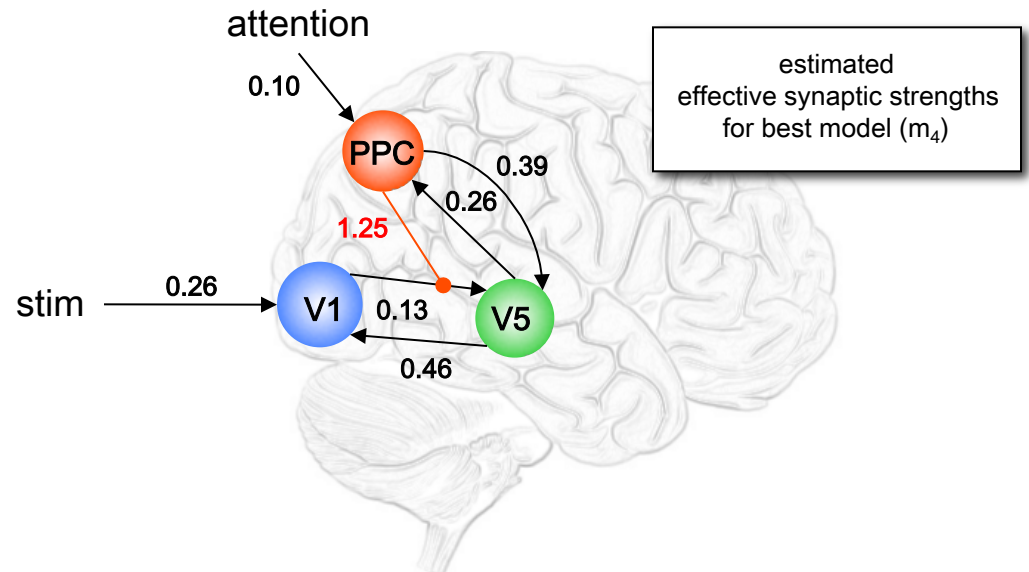
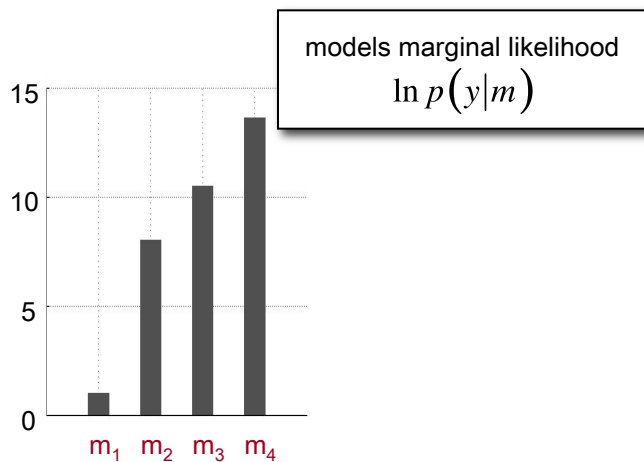
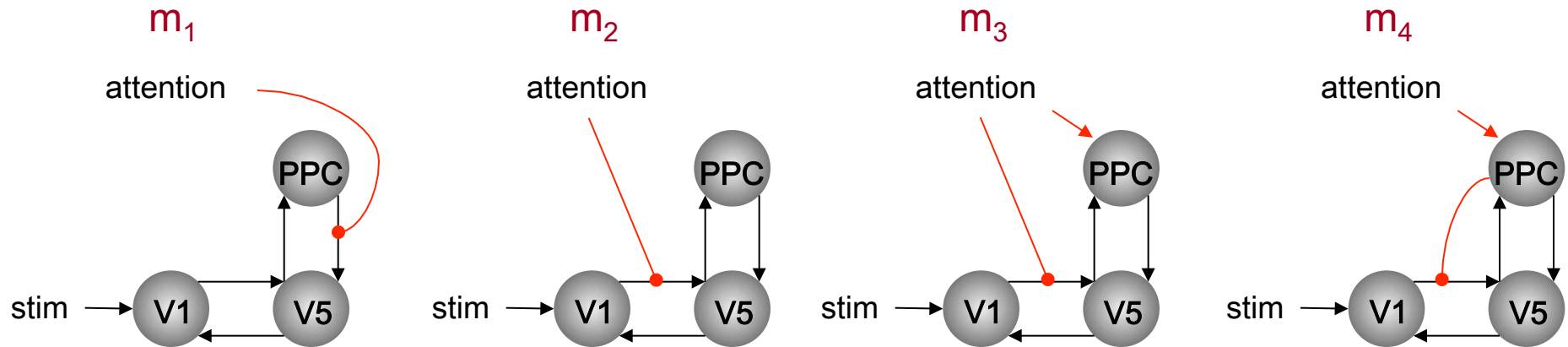


Neural model:
1 state variable per region
bilinear state equation
no propagation delays



Neural model:
8 state variables per region
nonlinear state equation
propagation delays

Bayesian Model Selection for fMRI



fMRI time series analysis with spatial priors

degree of smoothness

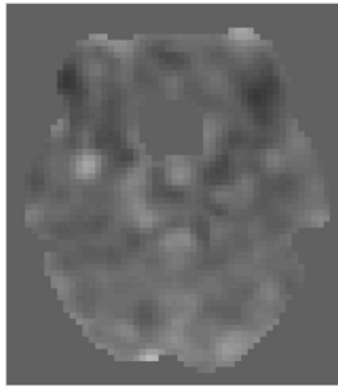
Spatial precision matrix

$$Y = X\beta + \varepsilon$$

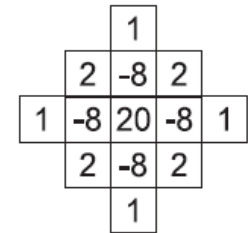
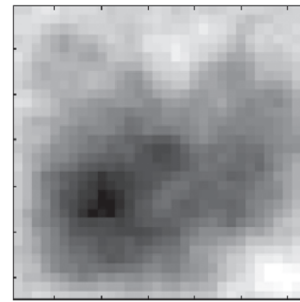
$$p(\beta) = N(0, \underline{\alpha}^{-1} \underline{L}^{-1})$$



aMRI



Smooth Y (RFT)

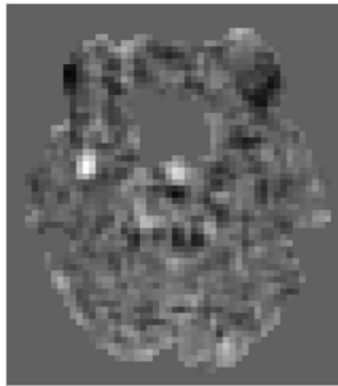


prior precision
of GLM coeff

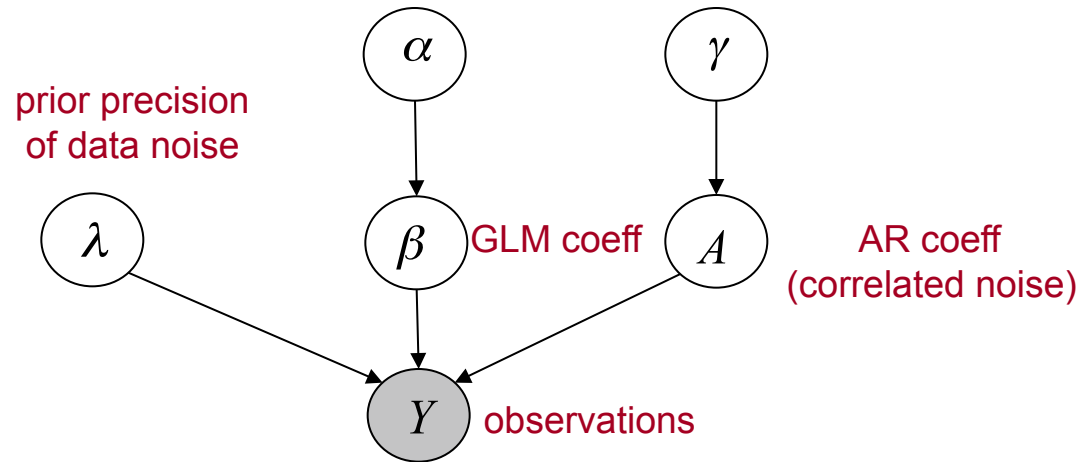
prior precision
of AR coeff



ML estimate of β



VB estimate of β



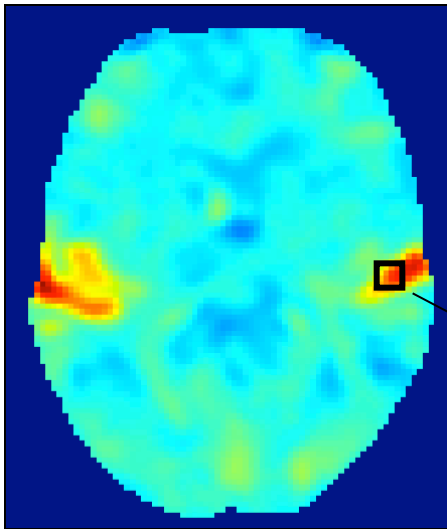
fMRI time series analysis with spatial priors:

posterior probability maps

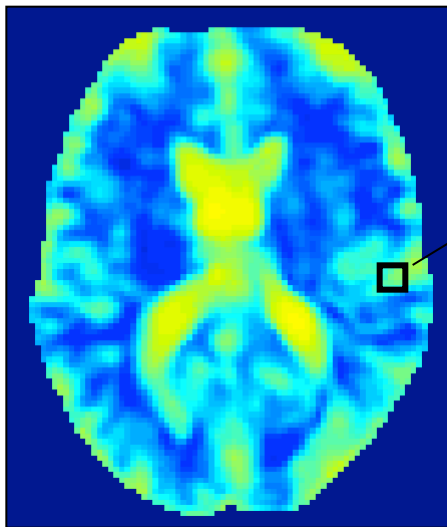
Display only voxels that exceed e.g. 95%

$$p > p_{th}$$

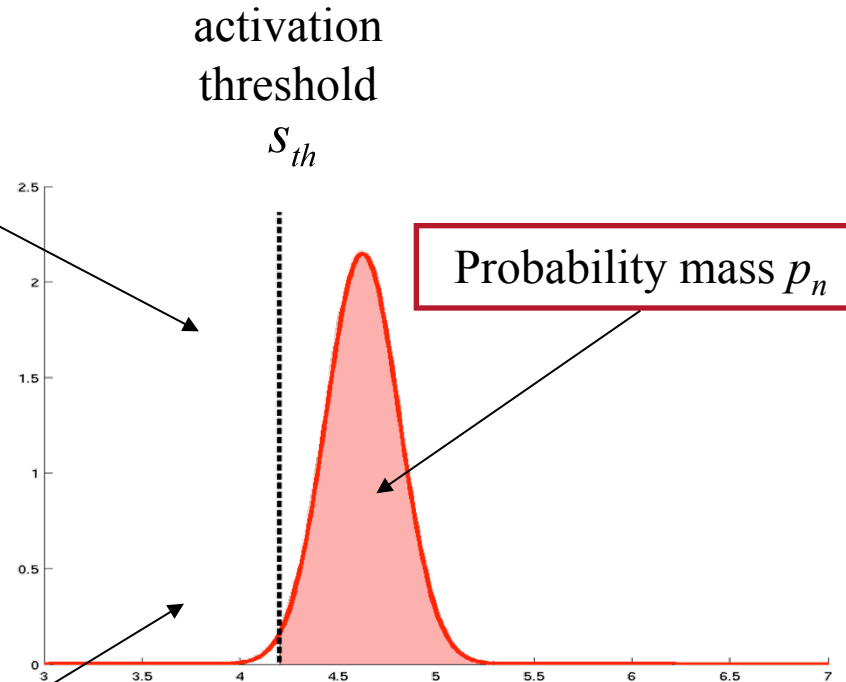
$$p = q(\beta > s_{th})$$



Mean (*Cbeta_*.img*)



Std dev (*SDbeta_*.img*)

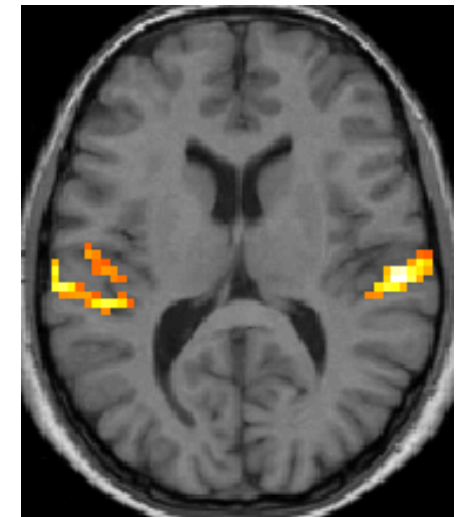


probability of getting an effect, given the data

$$q(\beta_n) = N(\mu_n, \Sigma_n)$$

mean: *size of effect*

covariance: *uncertainty*



PPM (*spmP_*.img*)

fMRI time series analysis with spatial priors: *Bayesian model selection*

$$\log p(y|m) \approx F(q)$$

Log-evidence maps

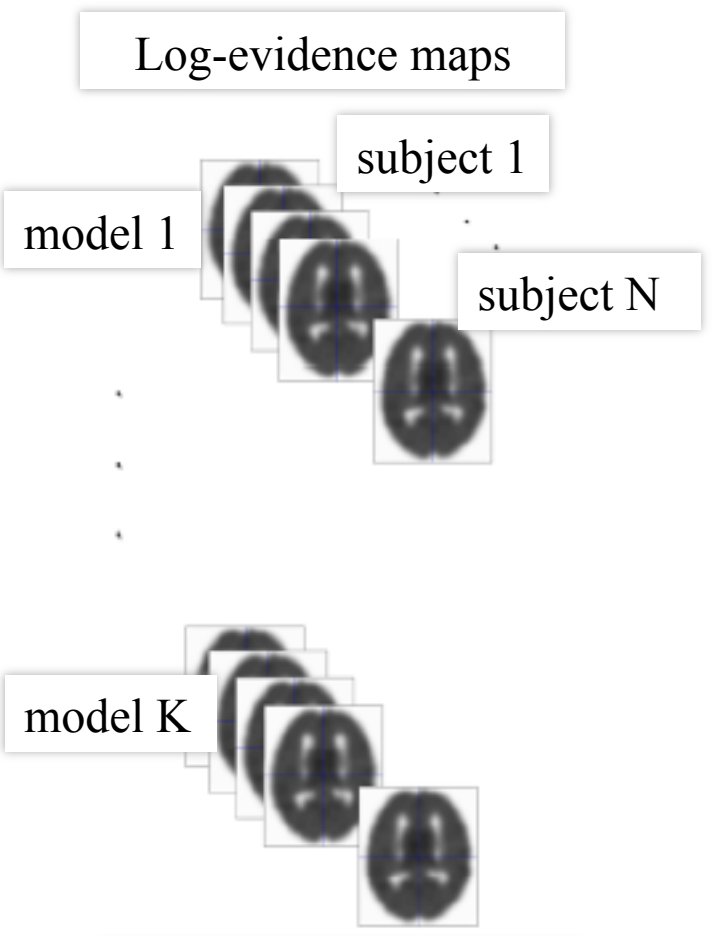
subject 1

model 1

subject N

model K

Compute log-evidence
for each model/subject



fMRI time series analysis with spatial priors: *Bayesian model selection*

$$\log p(y|m) \approx F(q)$$

Log-evidence maps

subject 1

model 1

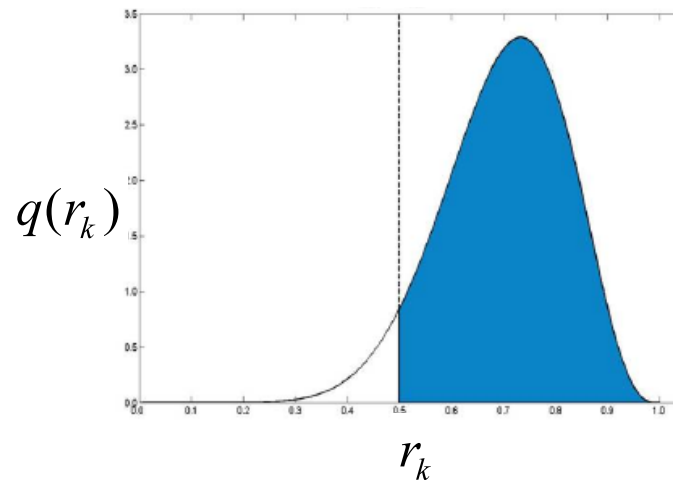
subject N

model K

Compute log-evidence
for each model/subject

BMS maps

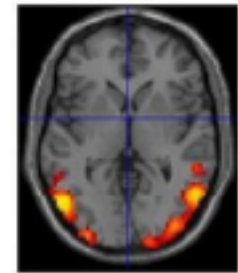
$$q(r_k > 0.5) = 0.941$$



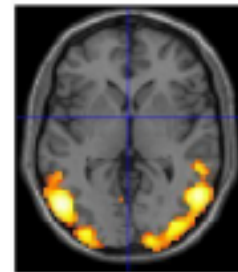
Probability that model k
generated data

$$\langle r_k \rangle > \gamma$$

$$\varphi_k > \gamma$$



PPM

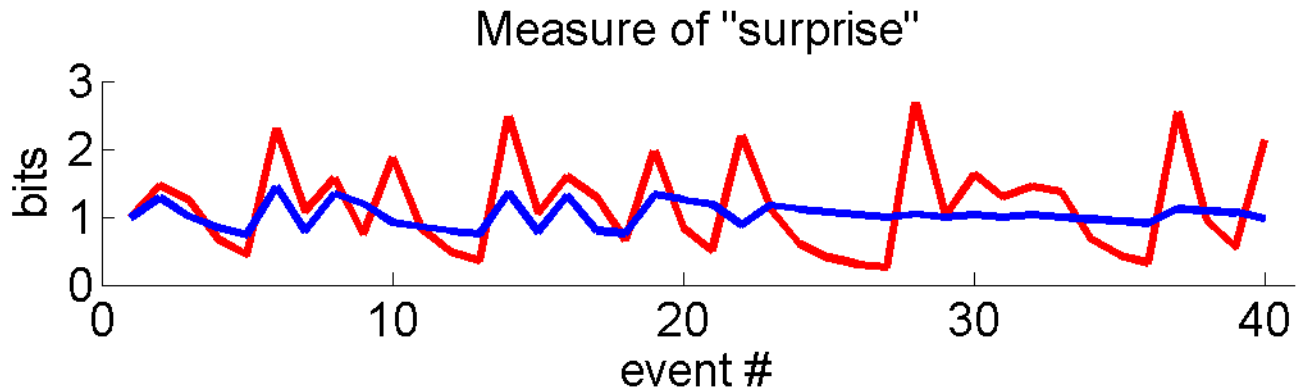
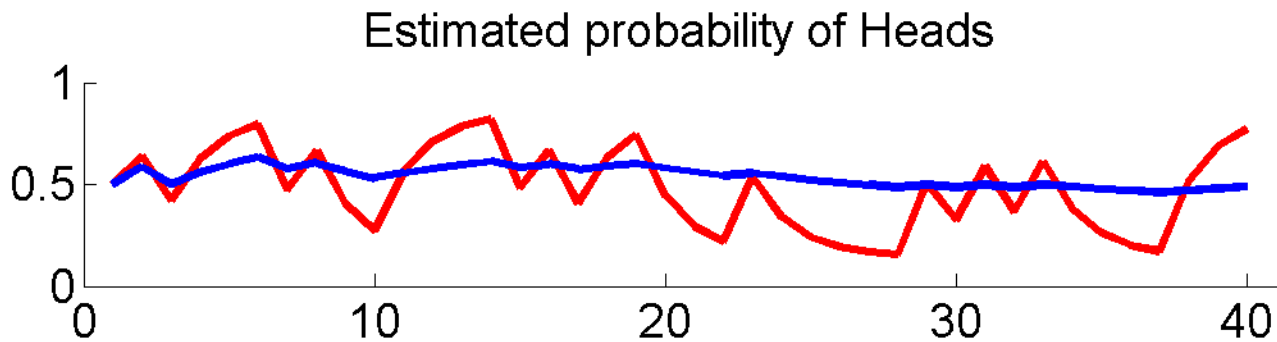
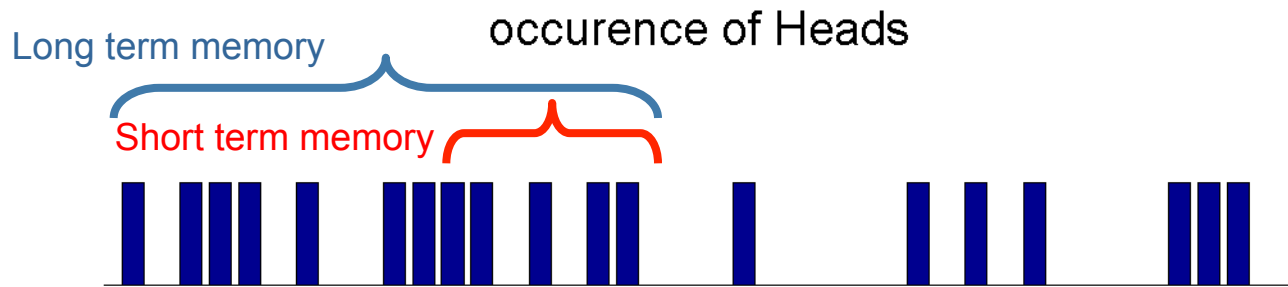


EPM

model k

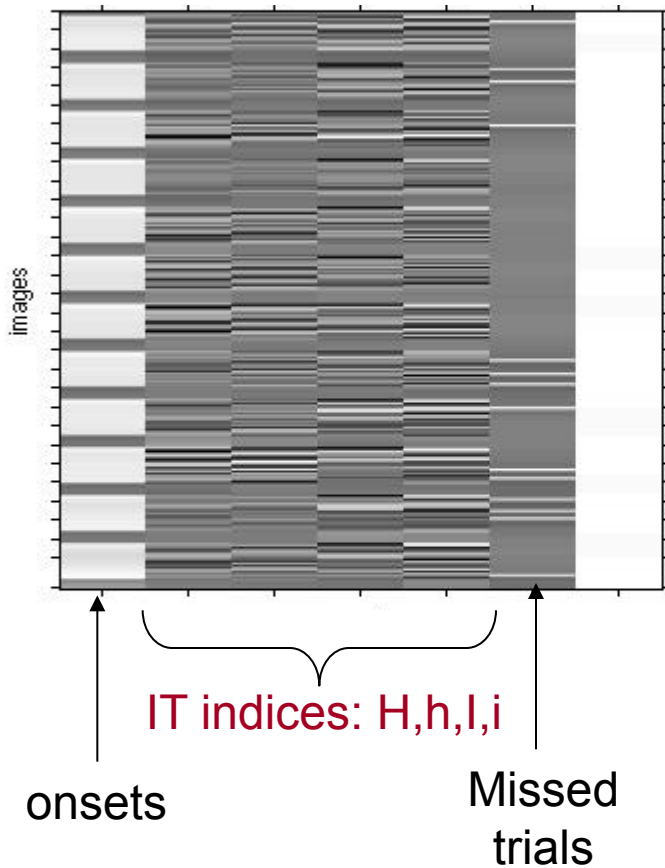
Joao et al, 2009

Reminder...

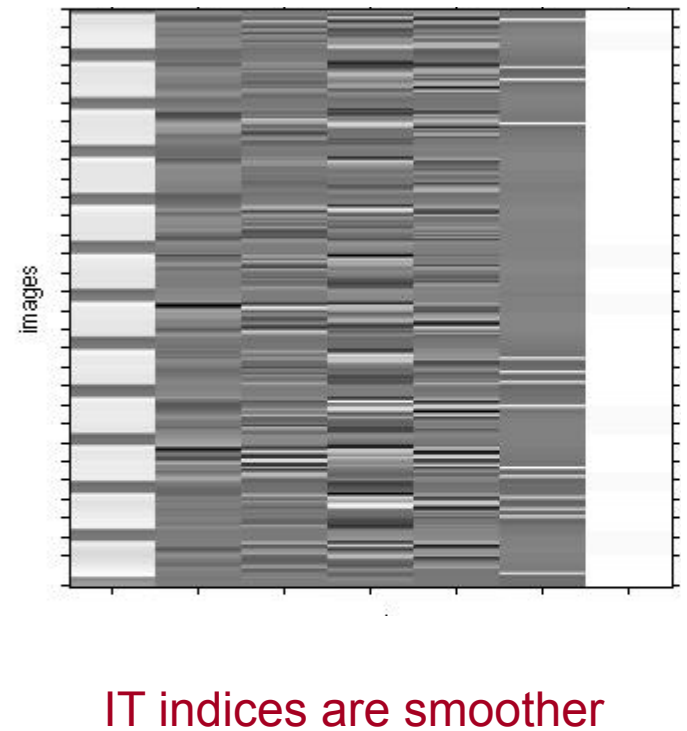


Compare two models

Short-term memory model



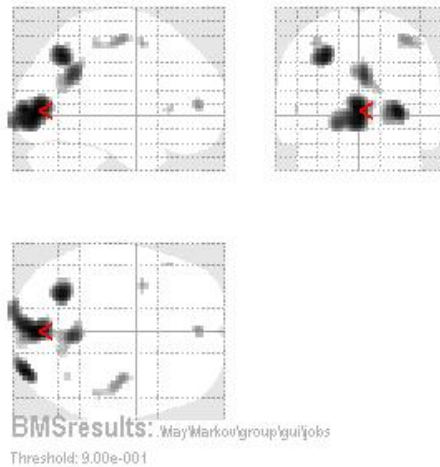
long-term memory model



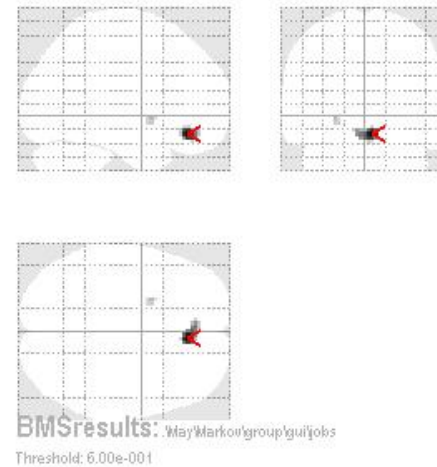
H=entropy; h=surprise; I=matural information; i=matural surprise

Group data: Bayesian Model Selection maps

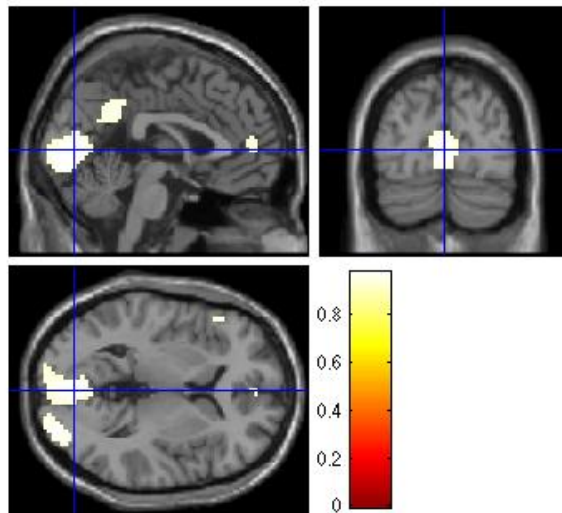
Regions best explained by short-term memory model



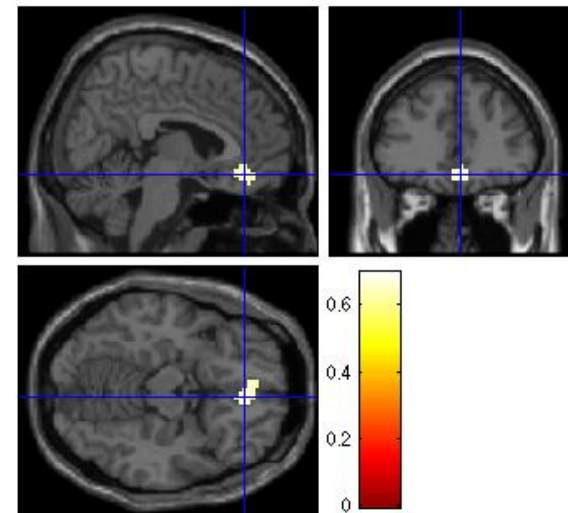
Regions best explained by long-term memory model



primary visual cortex



frontal cortex (executive control)



Thank-you