

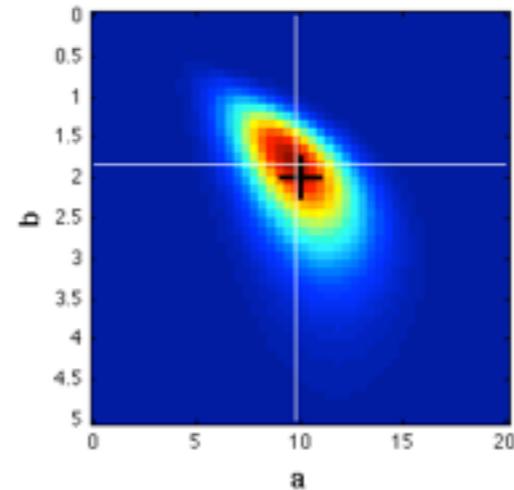
Bayesian modelling

Saad Jbabdi

- Conditioning and marginalisation
- Examples of Bayesian inference
- Approximate inference
 - Analytic
 - Sampling

$$p(A, B)$$

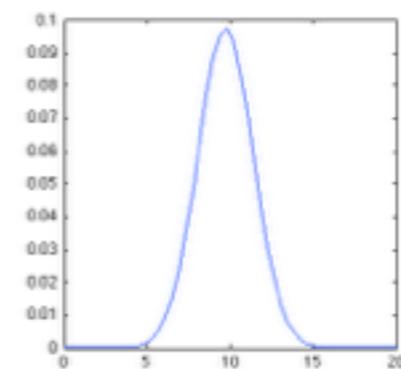
Joint distribution



$$\int p(A, B) dA dB = 1$$

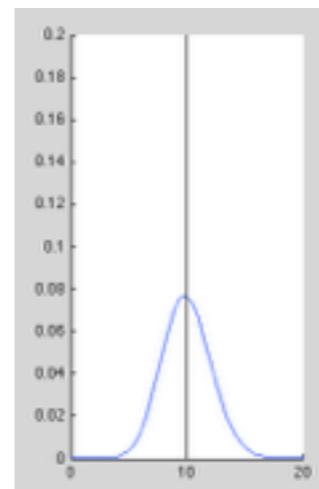
Conditioning

$$p(A | B)$$



Marginalisation

$$p(A, \cancel{B})$$



$$p(A) = \int p(A, B) dB$$

- Relationship between **joint probability** and **conditional probability**

$$p(A, B) = p(A | B) * p(B)$$
 product rule

- Relationship between **joint probability** and **marginal probability**

$$p(A) = \sum_B \{ p(A, B) \}$$
 sum rule

Bayes “theorem”

$$P(A, B) = P(A | B) * P(B)$$

product rule

$$P(A, B) = P(B | A) * P(A)$$

product rule



$$P(A | B) = P(B | A) * P(A) / P(B)$$

Bayesian modelling and inference

y : data

a : parameters of a model $y=F(a)$

$$\frac{p(a|y)}{\text{infer}} = \frac{p(y|a) * p(a)}{\text{model}} / p(y)$$

Bayesian modelling and inference

$$P(a | y) = P(y | a) * \text{prior} / P(y)$$

- model = prior and likelihood
- Bayesian inference = find the posterior

A simple example

$$y = a + \text{noise}$$

I measure one noisy data point

what is the probability distribution of a (given y)?

A simple example

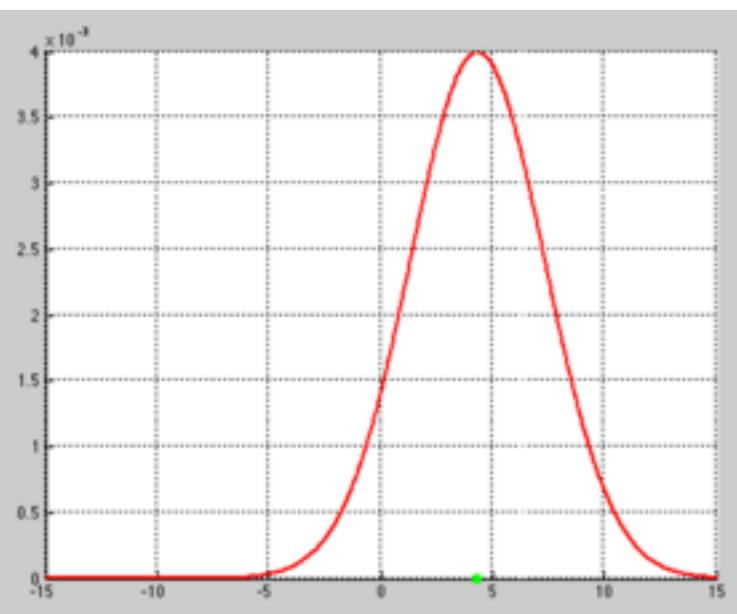
$$y = a + \text{noise}$$

assumption : noise is $N(0,s^2)$

now, we have a **likelihood**:

$$P(y|a) = N(y|a,s^2)$$

$$p(y|a) = (2\pi s^2)^{-1} \exp\left(-\frac{(y - a)^2}{2s^2}\right)$$



a

A simple example

$$y = a + \text{noise}$$

$$p(y|a) = N(y|a, s^2) \quad \text{likelihood}$$

$$p(a) = N(a|a_0, s_0^2) \quad \text{prior}$$

$$p(a|y) = N(y|a, s^2) \ N(a|a_0, s_0^2) / p(y)$$

(product of 2 Gaussians)

$$p(a|y) \sim \exp(-0.5(y-a)^2/s^2) * \exp(-0.5(a-a_0)^2/s_0^2)$$

$$= \exp(-0.5*(a-a_p)^2/s_p^2 + \text{const})$$

$$\sim N(a|a_p, s_p^2)$$

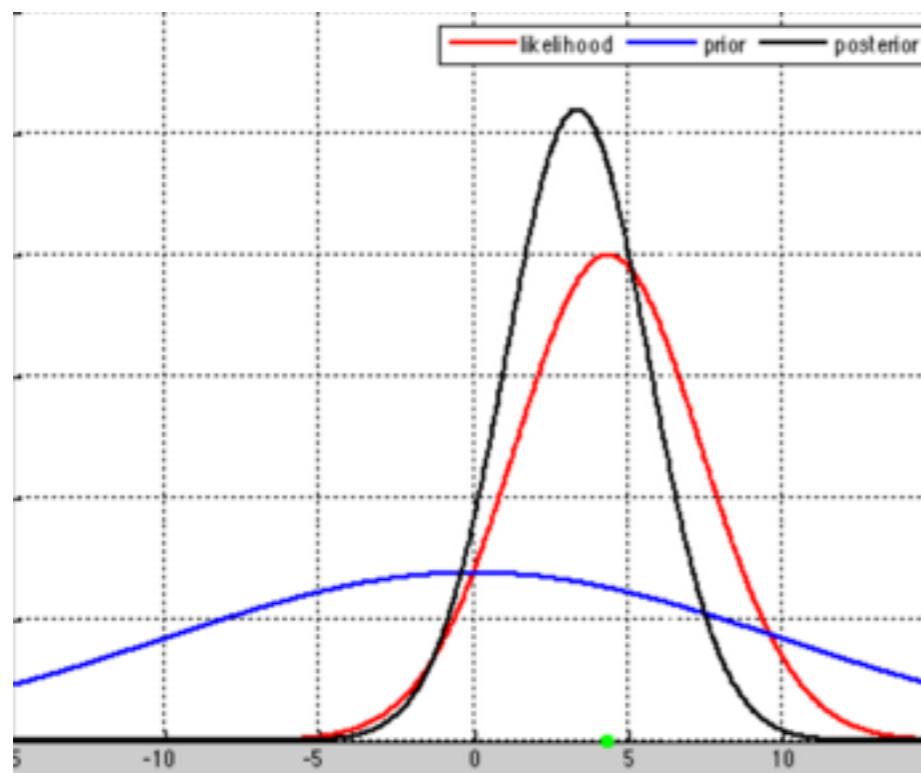
$$P(y|a) = N(y|a, s^2) \quad \text{likelihood}$$

$$P(a) = N(a|a_0, s_0^2) \quad \text{prior}$$

$$P(a|y) = N(a|a_p, s_p^2) \quad \text{posterior}$$

$$\begin{aligned} a_p &= s_p^2(\beta y + \beta_0 a_0) \\ s_p^2 &= (\beta + \beta_0)^{-1} \end{aligned}$$

precision (1/variance)



multiple data points?

extra assumption: noise is i.i.d.
(independent identically distributed samples)

$$P(y_1, y_2, \dots, y_n | a) = P(y_1 | a) P(y_2 | a) \dots P(y_n | a)$$

independent

$$P(y_1, y_2, \dots, y_n | a) = N(y_1 | a, s^2) N(y_2 | a, s^2) \dots N(y_n | a, s^2)$$

identically distributed

A simple example

$$y = a + \text{noise}$$

I measure N data points

what is the probability distribution of a?

$$\prod_{i=1}^N \exp\left(-\frac{(y_i - a)^2}{2s^2}\right) \rightarrow \exp\left(-\sum_{i=1}^N \frac{(y_i - a)^2}{2s^2}\right)$$

$$\begin{aligned} p(a|y) \sim p(y|a)*p(a) &= \exp(-0.5\sum(y_i-a)^2/s^2)*\exp(-0.5(a-a_0)^2/s_0^2) \\ &= \exp(-0.5*(a-a_p)^2/s_p^2+\text{const}) \\ &\sim N(a|a_p, s_p^2) \end{aligned}$$

A simple example

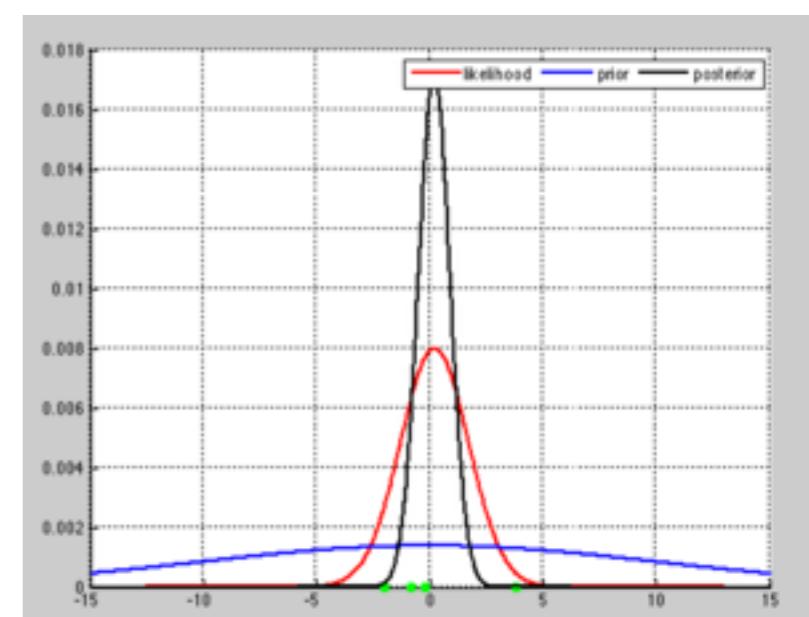
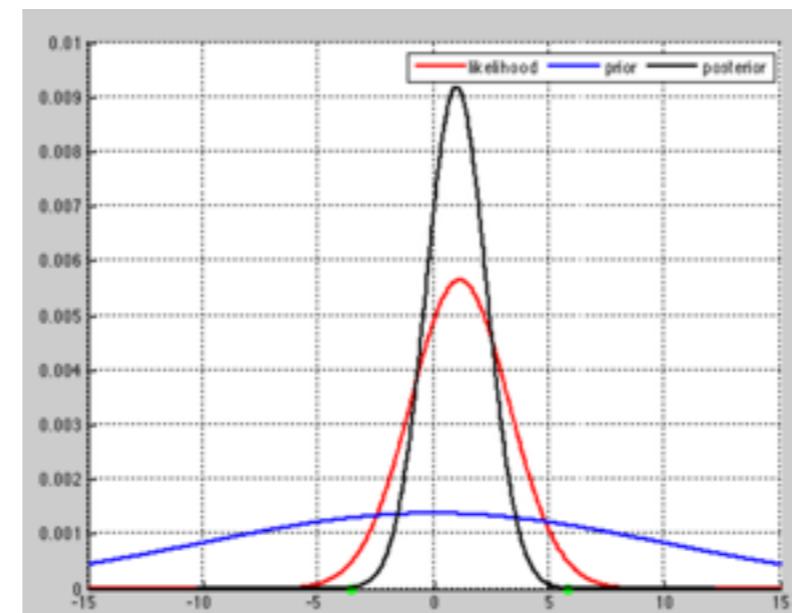
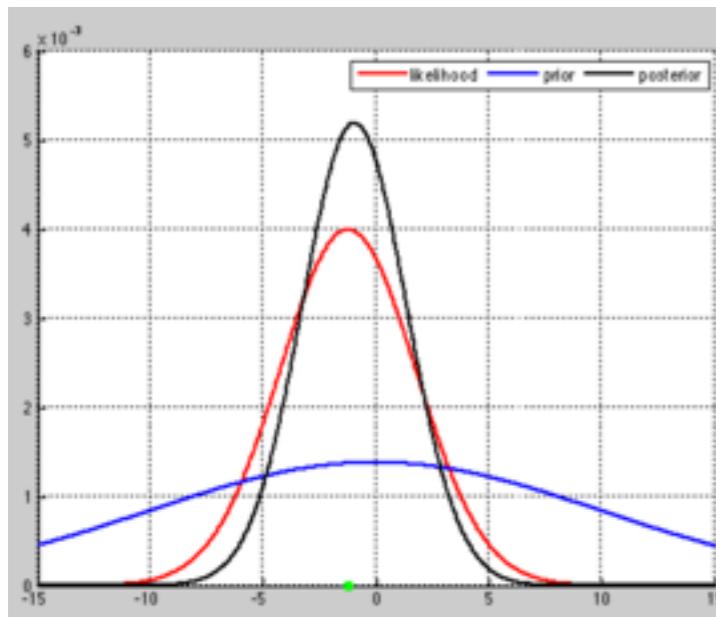
$$y = a + \text{noise}$$

$$p(a|y) \propto N(a_p, s_p^2)$$

$$a_p = s_p^2 * \left(n\beta \frac{\sum y_i}{n} + \beta_0 a_0 \right)$$

$$s_p^2 = (n\beta + \beta_0)^{-1}$$

precision (1/variance)



Regression

$$y = a * x$$

I measure one noisy data point

what is the probability distribution of a?

Regression

$$y = a^*x$$

noise model: $y = a^*x + N(0, s^2)$

likelihood: $p(y | a) = N(a^*x, s^2)$

prior: $p(a) = N(a|a_0, s_0^2)$

Bayes: $p(a | y) \sim p(y | a) * p(a)$

$$\begin{aligned} p(a|y) &\sim p(y|a) * p(a) = \exp(-0.5 \sum (y_i - ax_i)^2 / s^2) * \exp(-0.5(a - a_0)^2 / s_0^2) \\ &= \exp(-0.5 * (a - a_p)^2 / s_p^2 + \text{const}) \\ &\sim N(a | a_p, s_p^2) \end{aligned}$$

Regression

$$p(a|y) \propto N(a_p, s_p)$$

$$a_p =$$

$$\frac{bx^T y + b_0 a_0}{bx^T x + b_0}$$

$$s_p =$$

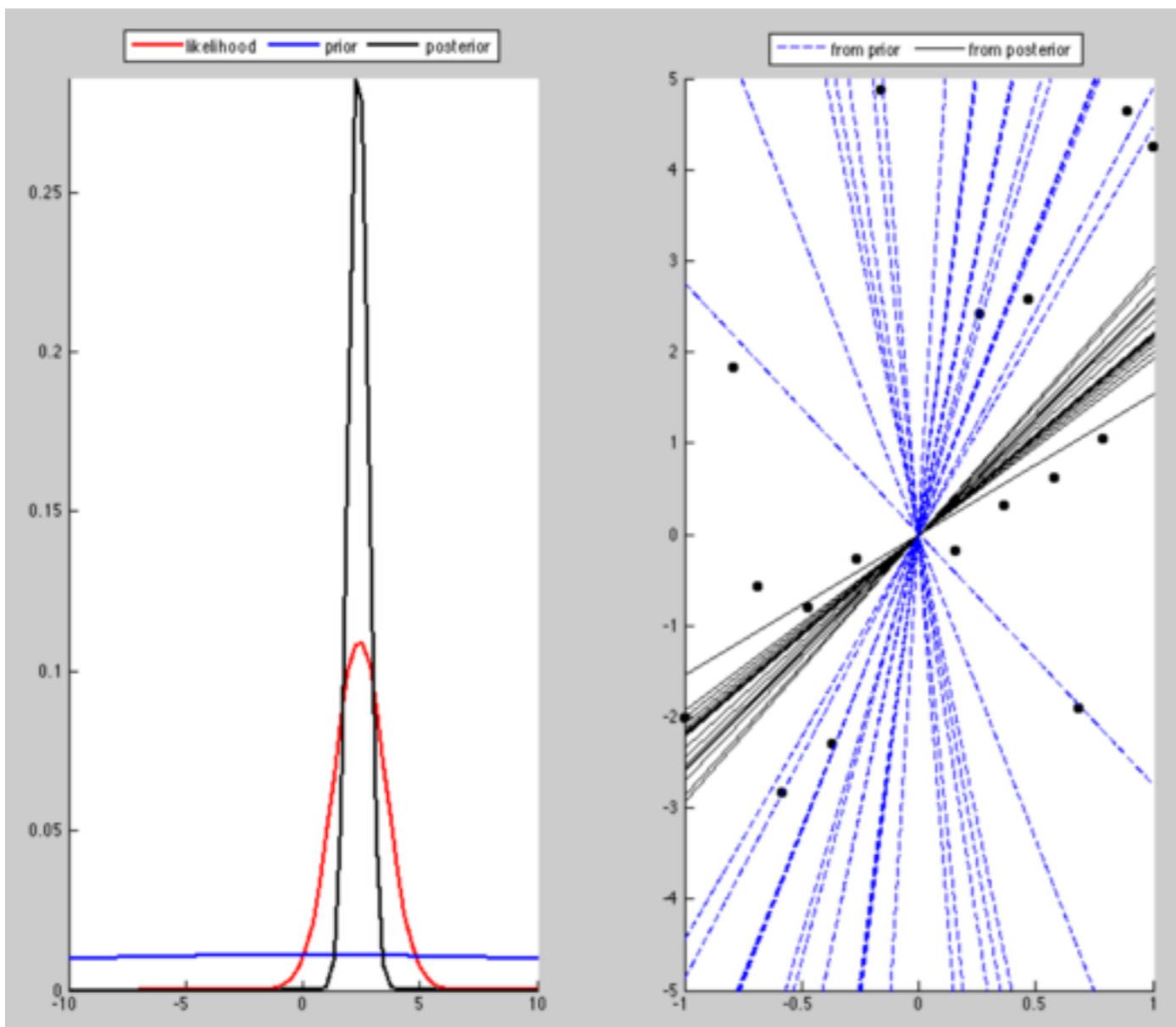
$$(bx^T x + b_0)^{-1}$$

$$b = 1/s^2$$

$$b_0 = 1/s_0^2$$

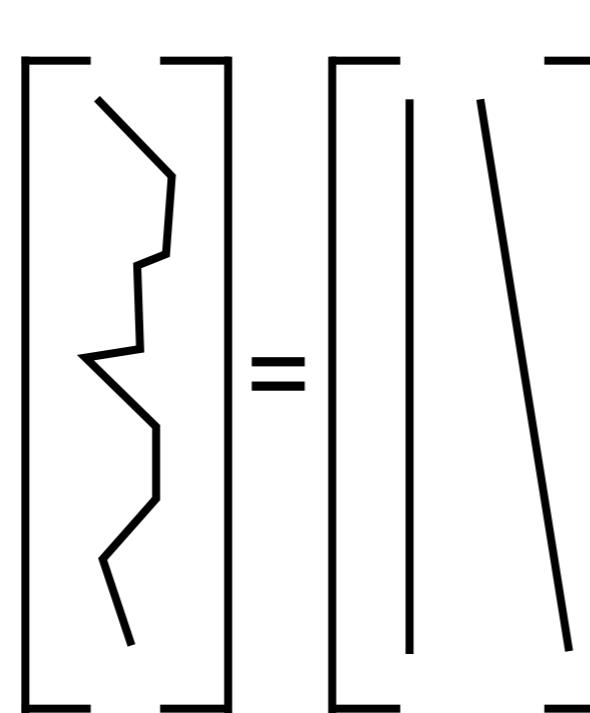
$$(X^T X)^{-1} X^T Y$$

least squares regression!



Sequential learning

$$Y = X * b$$



The diagram shows a sequence of vertical segments forming a zigzag pattern, representing a learned function or sequence. This is followed by an equals sign (=). To the right of the equals sign is a multiplication symbol (*). To the right of the multiplication symbol is a vector [b₁ b₂] with two arrows pointing to it. The top arrow is labeled "intercept" and the bottom arrow is labeled "slope".

Sequential learning

$$P(Y | b) = P(y_1|b)*P(y_2|b)*...*P(y_n|b)$$

i.i.d. noise

$$P(y_i|b)=N(y_i|(Xb)_i, s^2)$$

$$P(b | Y) = P(y_n|b)*P(y_{n-1}|b)*...*P(y_2|b)*P(y_1|b) * p(b)$$



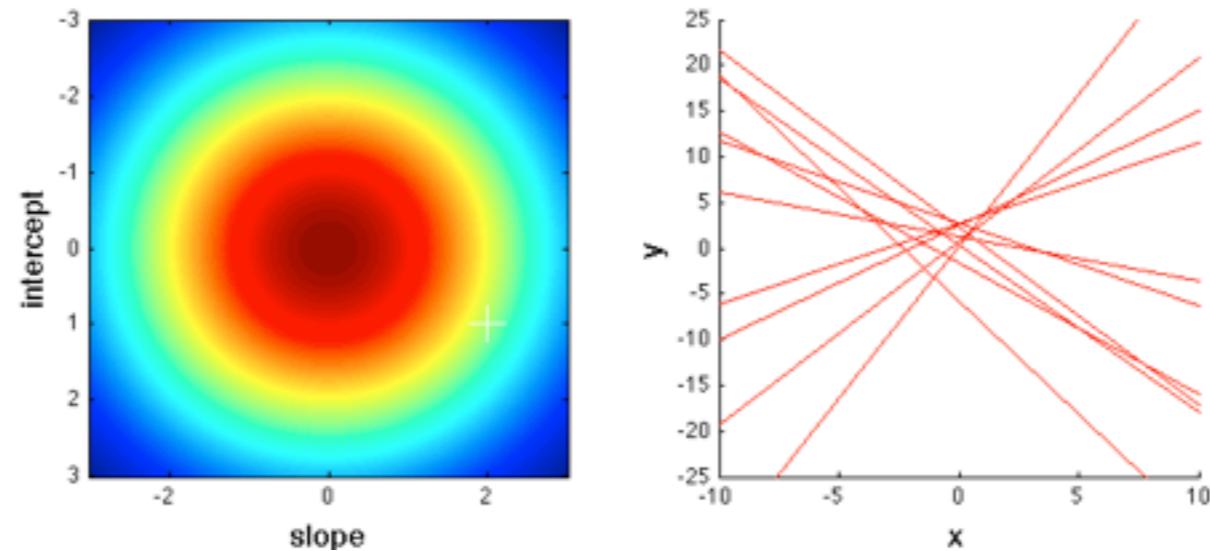
$$P(b|y_1)$$



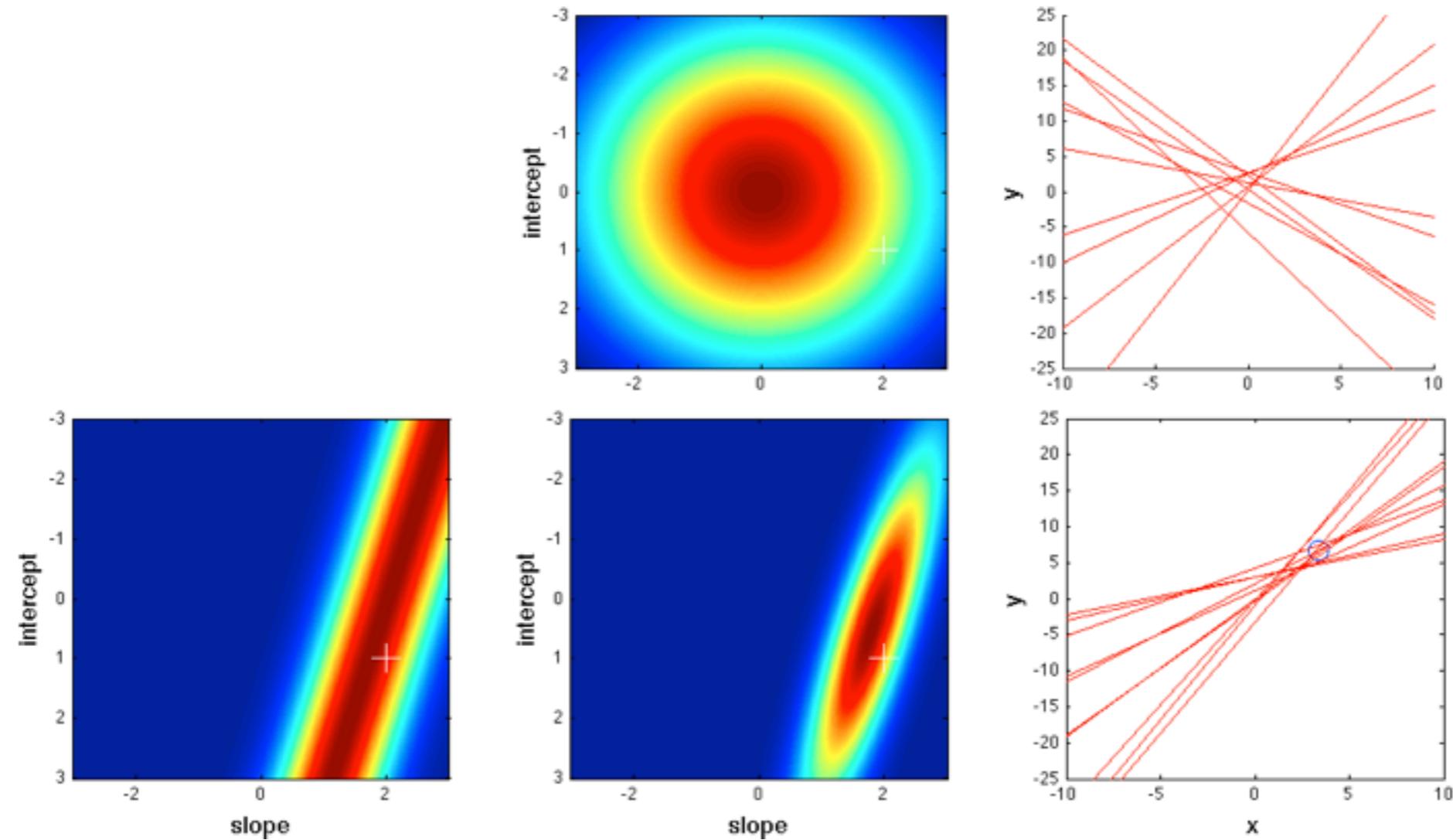
$$P(b|y_1, y_2)$$

posterior given y_1 becomes prior for y_2 , etc.etc.

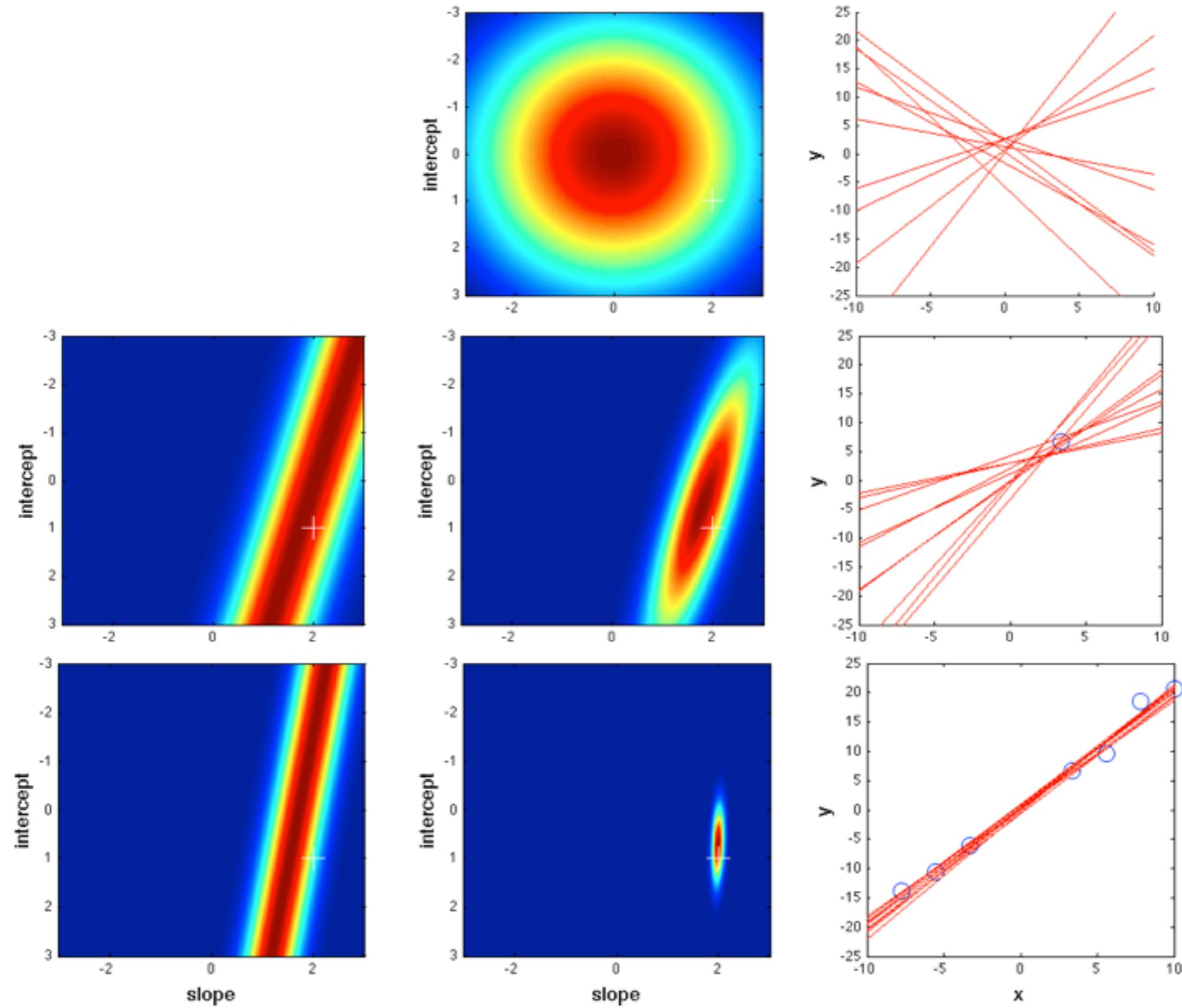
Sequential learning



Sequential learning



Sequential learning



Sequential learning

The posterior distribution is the prior for the next
data point

Nonlinear model

$$y = a * \exp(-b * x) = f(a, b)$$

noise model: $y = f(a, b) + N(0, s^2)$

likelihood: $p(y|a,b) = N(y|f(a,b), s^2)$

prior: $p(a,b) = N(a|a_0, s_{a0}^2) * N(b|b_0, s_{b0}^2)$

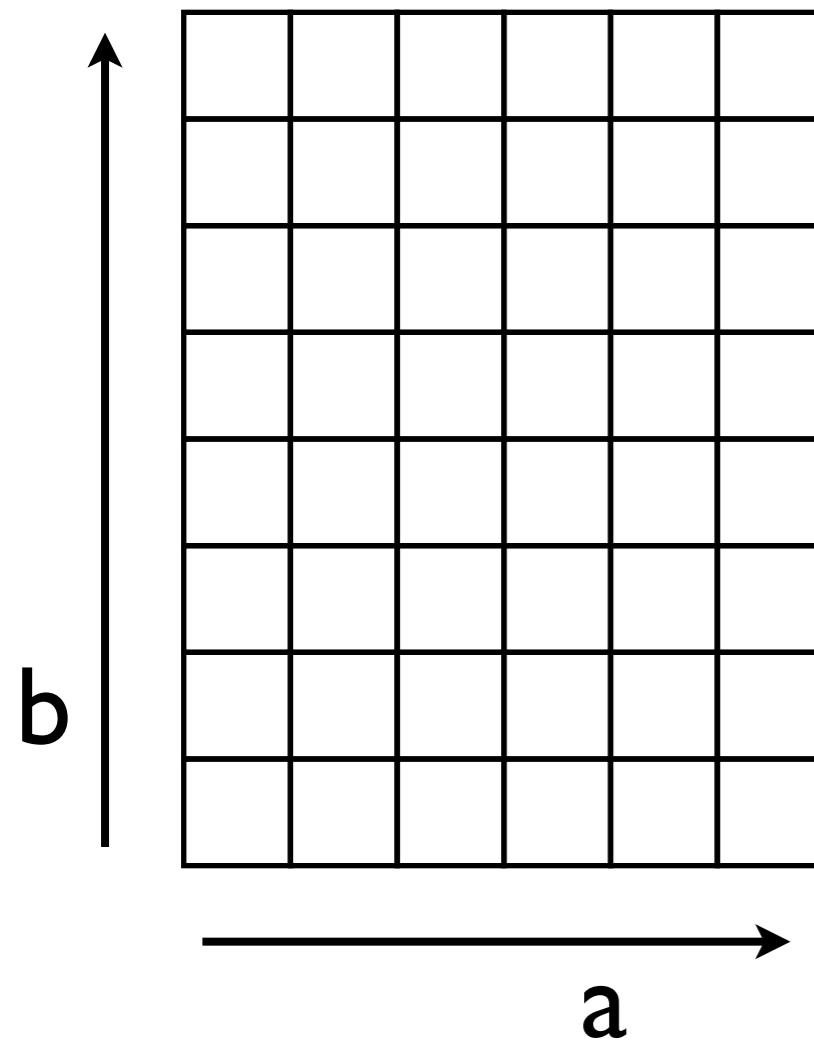
posterior: $p(a,b|y) \sim \exp[-0.5(y-f(a,b))^2/s^2] * p(a,b)$

not a Gaussian... what do we do?

Grid evaluation

$$y = a * \exp(-b * x)$$

posterior: $p(a,b|y) = \text{Const} * \exp[-0.5(y-f(a,b))^2/s^2] * p(a,b)$



calculate $p(a,b|y)$ (to a constant)
normalise by sum over the whole grid

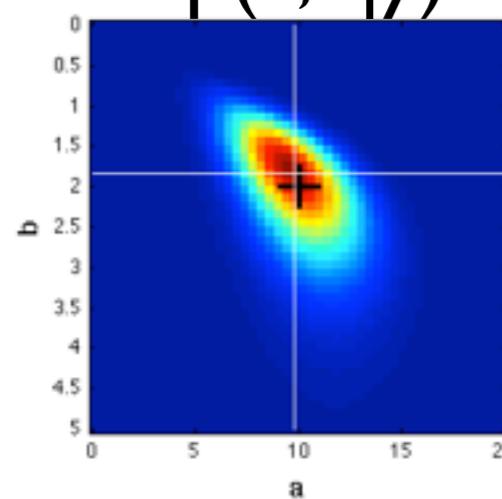
Grid evaluation

$$y = a * \exp(-b * x)$$

$$x = 1:10$$

full posterior
 $p(a,b|y)$

$a=10$
 $b=2$
 $s^2=3$

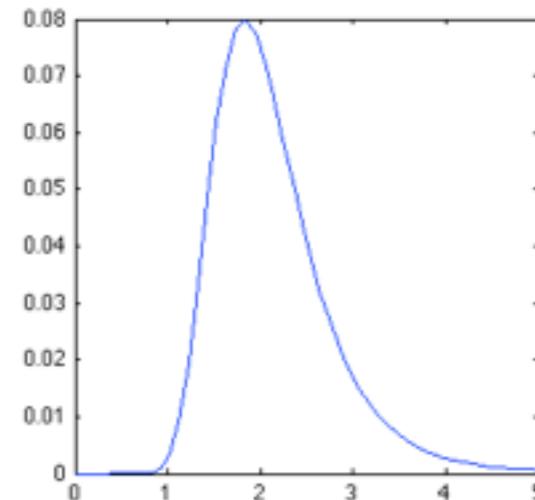


$p(a|y, b=2)$

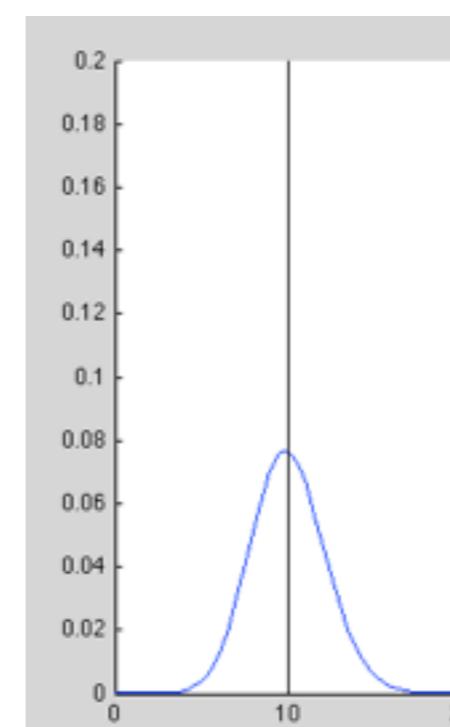
$p(a|y)$

$p(b|y)$

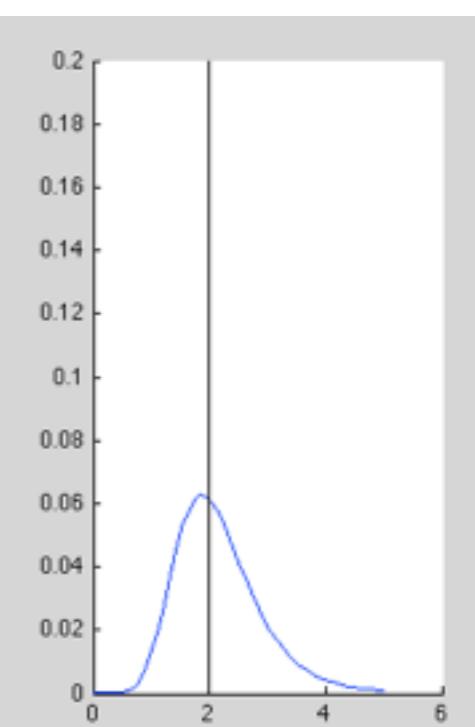
conditionals



$p(b|y, a=10)$



marginal posteriors



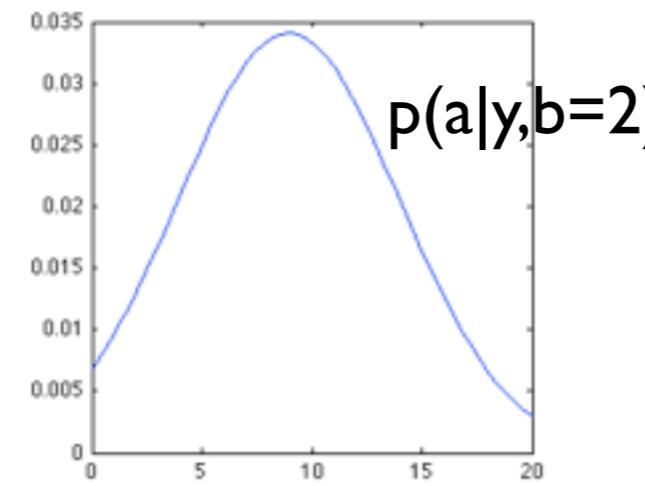
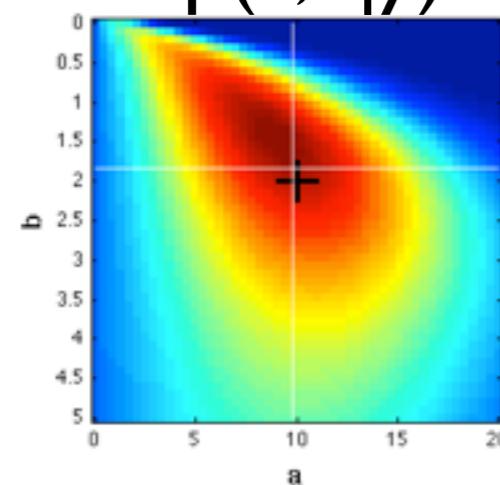
Grid evaluation

$$y = a * \exp(-b * x)$$

$$x = 1:10$$

full posterior
 $p(a,b|y)$

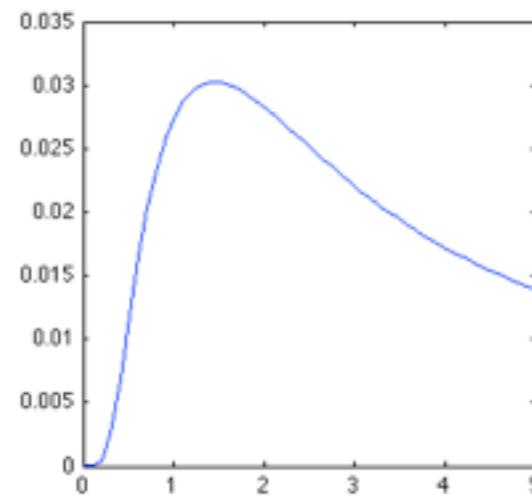
$a=10$
 $b=2$
 $s^2=10$



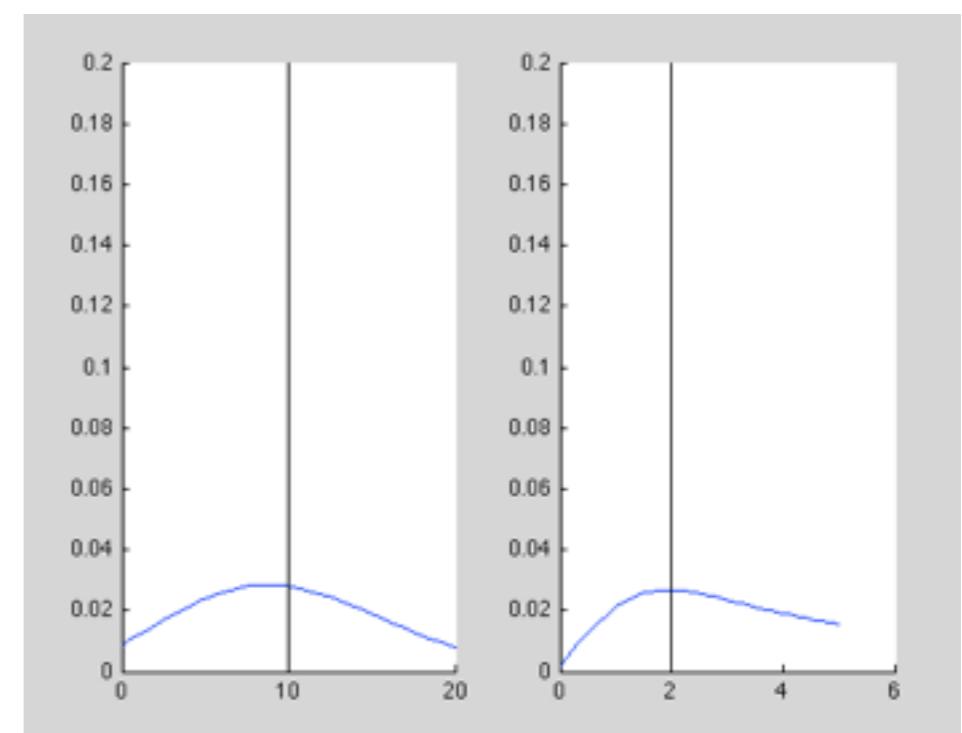
$p(a|y)$

$p(b|y)$

conditionals



$p(b|y, a=10)$



marginal posteriors

Grid search

but what if we had MANY parameters?

e.g. 5 parameters, 100 bins each = 100^5 values
to calculate (and store)

solutions:

- approximate with an easier distribution
- sample directly from the distribution (e.g. MCMC)

Laplace Approximation

posterior: $p(x|y) \sim \exp(-.5 * (y-f(x))^2/s^2) * p(x)$

$-\log p(x|y) \sim (y-f(x))^2/s^2$ ignore priors

$-\log p(x|y) \sim g(x)/s^2$

$g(x) \sim g(x_m) + g'(x_m)(x-x_m) + .5*g''(x_m)(x-x_m)^2 + \dots$ (Taylor theorem)

$g(x) \sim g(x_m) + .5*g''(x_m)(x-x_m)^2$

constant



quadratic, i.e. $p(x)$
is like a Gaussian

$x_m = \text{maximum likelihood}$

Laplace Approximation

$$\log(p(a,b|y)) \sim g(x)/s^2 \sim \text{const} + 0.5*g''(x_m)(x-x_m)^2/s^2$$

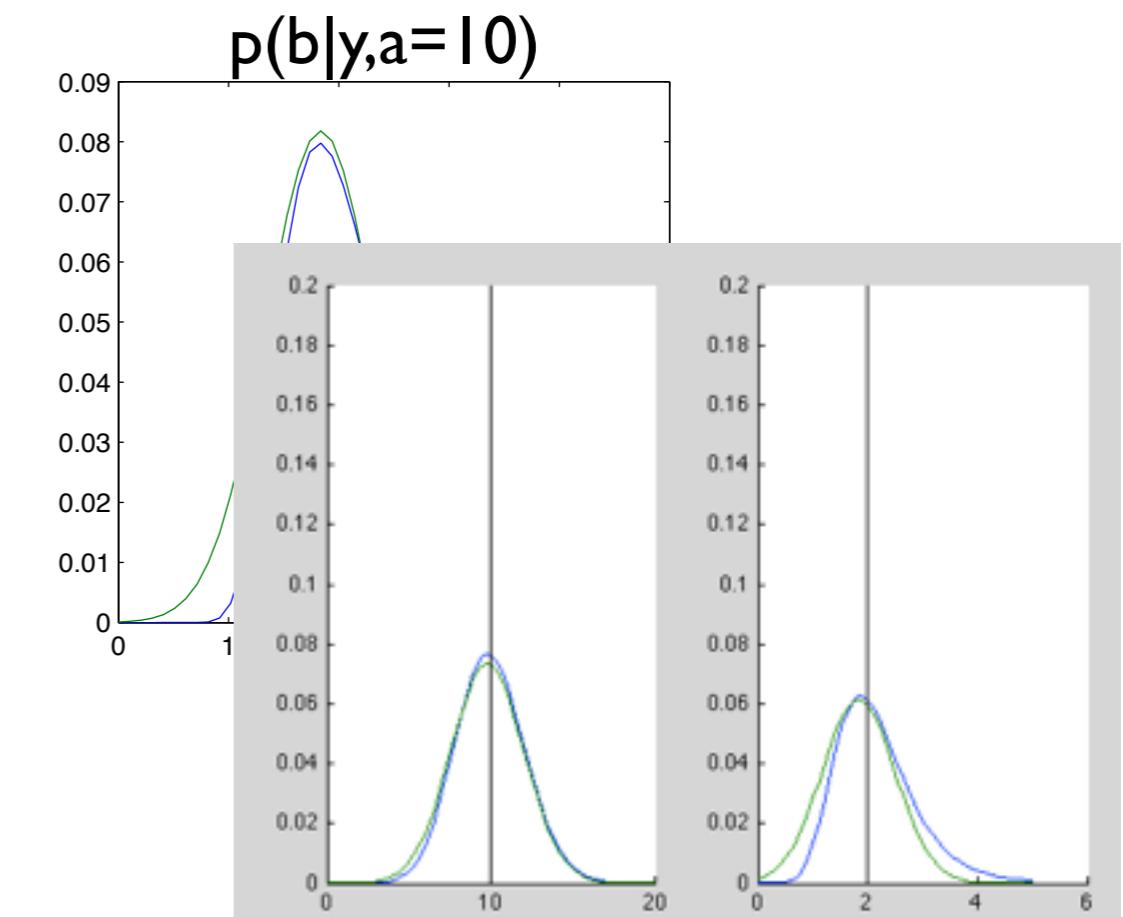
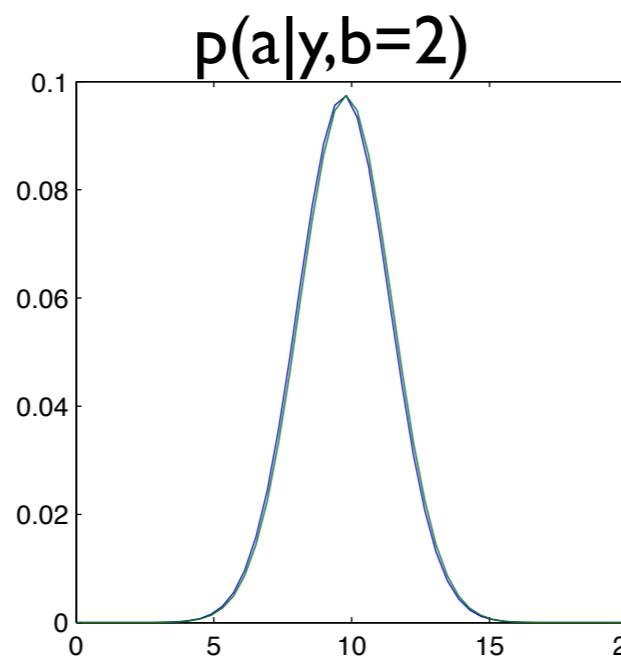
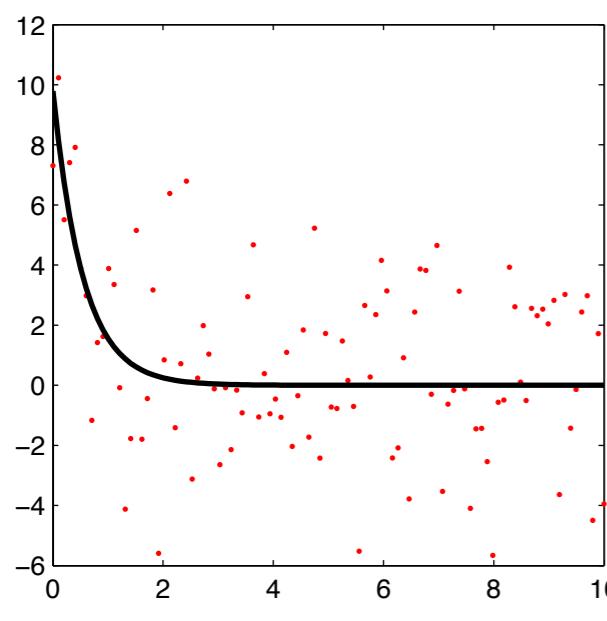
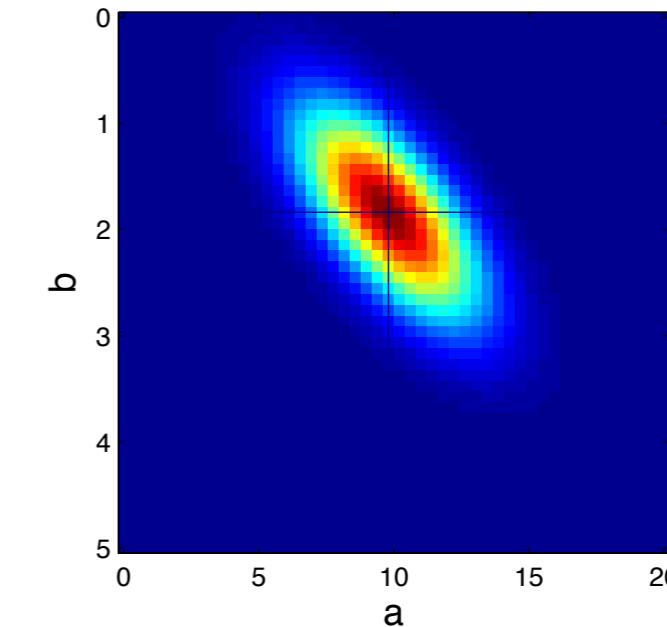
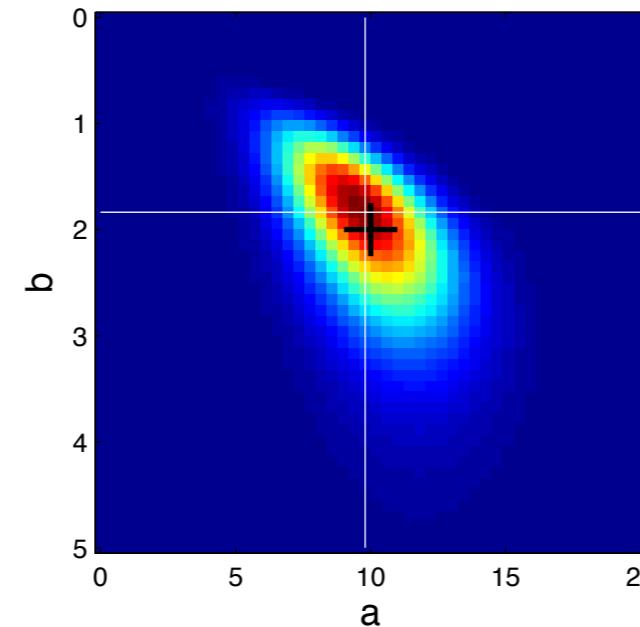
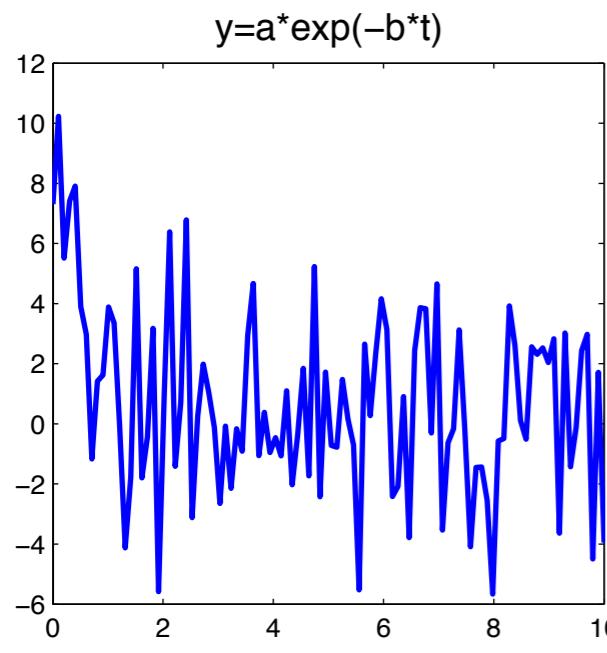
posterior: $p(a,b|y) \sim N(y | (a_m; b_m), s^2 H^{-1})$

$$H = \begin{bmatrix} g_{aa} & g_{ab} \\ g_{ba} & g_{bb} \end{bmatrix}$$

$$g = (y - f(a,b))^2$$

Laplace approx in one sentence:
“curvature of the posterior PDF
around its max”

Laplace Approximation



approximate full posterior by a “local” Gaussian
All marginals and conditionals are Gaussians

Cannot work if posterior is complicated

Variational Bayes (VB)

$q(\mathbf{a})$ instead of $p(\mathbf{a}|y)$

approximate with simple $q(\mathbf{a})$ with
unknown “parameters”

$$KL = \int q(\mathbf{a}) \log [p(\mathbf{a}|y)/q(\mathbf{a})]$$

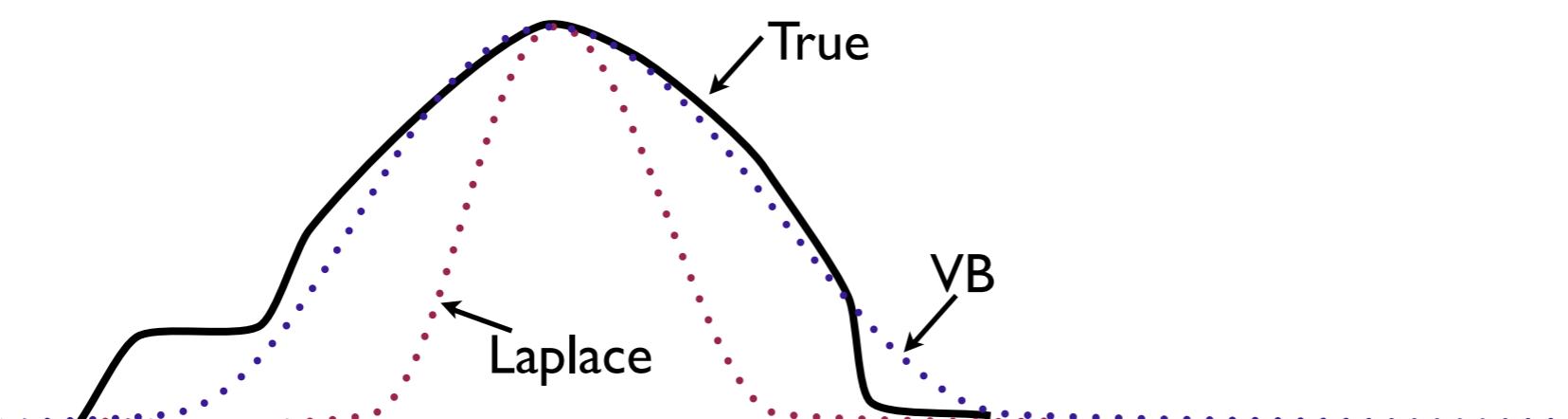
Kullback Leibler
(distance between true and
approximate posterior)

$$F = \int q(\mathbf{a}) \log [p(y|\mathbf{a})p(\mathbf{a})/q(\mathbf{a})]$$

define Free Energy

$$F+KL=\log[p(y)]$$

constant \rightarrow minimize KL wrt q by
maximizing F



The kind of maths that you typically have to do...

$$\begin{aligned}\ln q_\mu^*(\mu) &= \mathbb{E}_\tau[\ln p(\mathbf{X}|\mu, \tau) + \ln p(\mu|\tau) + \ln p(\tau)] + C \\&= \mathbb{E}_\tau[\ln p(\mathbf{X}|\mu, \tau)] + \mathbb{E}_\tau[\ln p(\mu|\tau)] + \mathbb{E}_\tau[\ln p(\tau)] + C \\&= \mathbb{E}_\tau[\ln \prod_{n=1}^N \mathcal{N}(x_n|\mu, \tau^{-1})] + \mathbb{E}_\tau[\ln \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})] + C_2 \\&= \mathbb{E}_\tau[\ln \prod_{n=1}^N \sqrt{\frac{\tau}{2\pi}} e^{-\frac{(x_n-\mu)^2\tau}{2}}] + \mathbb{E}_\tau[\ln \sqrt{\frac{\lambda_0\tau}{2\pi}} e^{-\frac{(\mu-\mu_0)^2\lambda_0\tau}{2}}] + C_2 \\&= \mathbb{E}_\tau\left[\sum_{n=1}^N \left(\frac{1}{2}(\ln \tau - \ln 2\pi) - \frac{(x_n - \mu)^2\tau}{2}\right)\right] + \mathbb{E}_\tau\left[\frac{1}{2}(\ln \lambda_0 + \ln \tau - \ln 2\pi) - \frac{(\mu - \mu_0)^2\lambda_0\tau}{2}\right] + C_2 \\&= \mathbb{E}_\tau\left[\sum_{n=1}^N -\frac{(x_n - \mu)^2\tau}{2}\right] + \mathbb{E}_\tau\left[-\frac{(\mu - \mu_0)^2\lambda_0\tau}{2}\right] + \mathbb{E}_\tau\left[\sum_{n=1}^N \frac{1}{2}(\ln \tau - \ln 2\pi)\right] + \mathbb{E}_\tau\left[\frac{1}{2}(\ln \lambda_0 + \ln \tau - \ln 2\pi)\right] + C_2 \\&= \mathbb{E}_\tau\left[\sum_{n=1}^N -\frac{(x_n - \mu)^2\tau}{2}\right] + \mathbb{E}_\tau\left[-\frac{(\mu - \mu_0)^2\lambda_0\tau}{2}\right] + C_3 \\&= -\frac{\mathbb{E}_\tau[\tau]}{2} \left\{ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right\} + C_3\end{aligned}$$

Noise trick

- so far we have assumed that we know the variance of the noise s^2
- what to do when we don't want to make that assumption?
 - Estimate it like the other parameters
 - Marginalise wrt the noise variance

Noise trick

Let us assume a Gaussian noise model.

$$p(y|x, s^2) = N(f(x), s^2)$$

If we want to estimate s^2 , we need a prior on s^2

Get one from a hat: $p(s^2) = 1/s^2$

$$p(x|y) = \int p(x|y, s^2) p(s^2) ds^2 \quad \text{Marginalisation}$$

whiteboard calculation of:

$$P(x|y) = \int N(y|f(x), s^2) P(s^2) ds^2$$

gives:

$$P(x|y) \propto |y - f(x)|^{-n}$$

- When we don't care about a parameter, but we still know it's “there”, it is often useful to “integrate it out” (marginalise)

likelihood

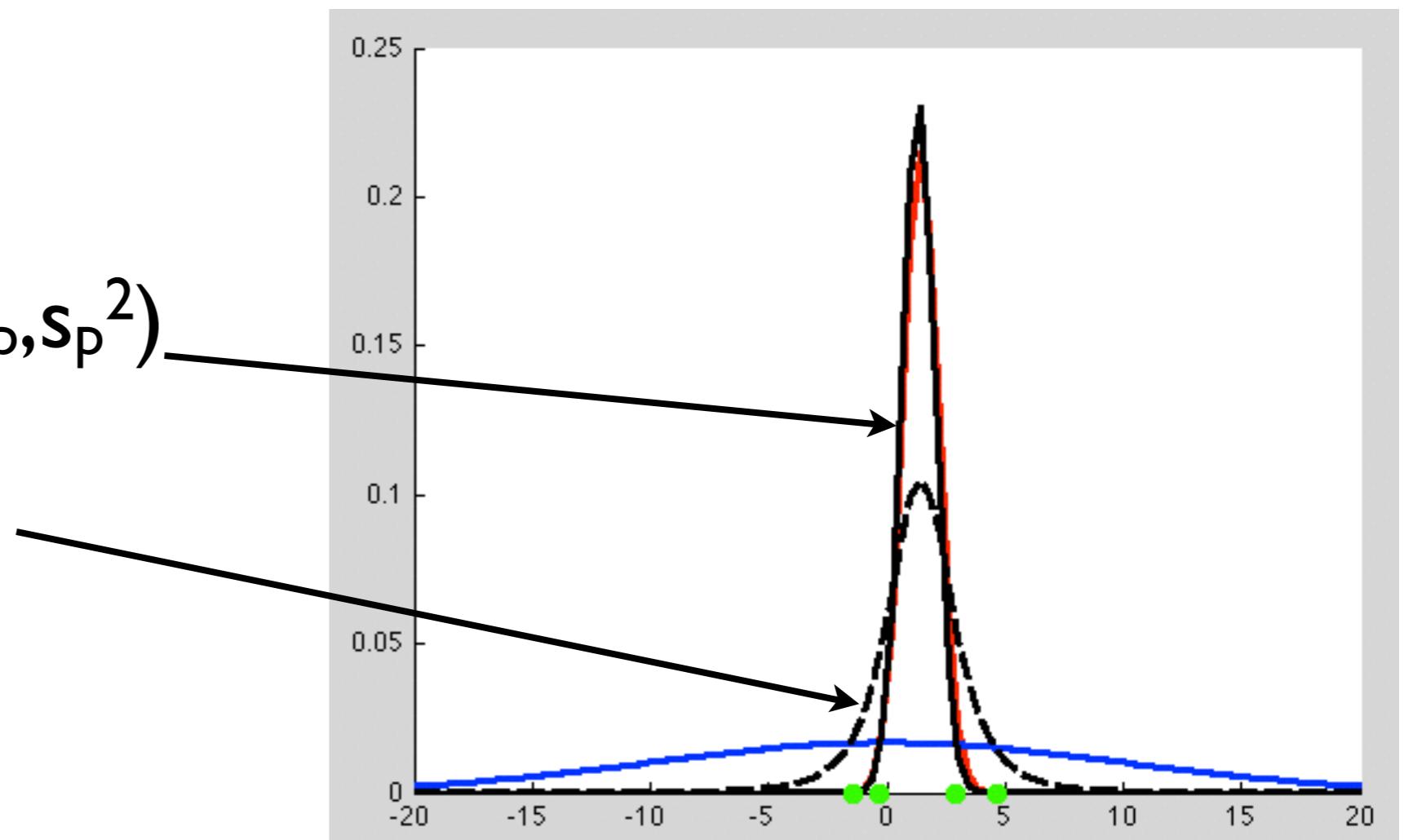
prior

$y = a + \text{noise}$

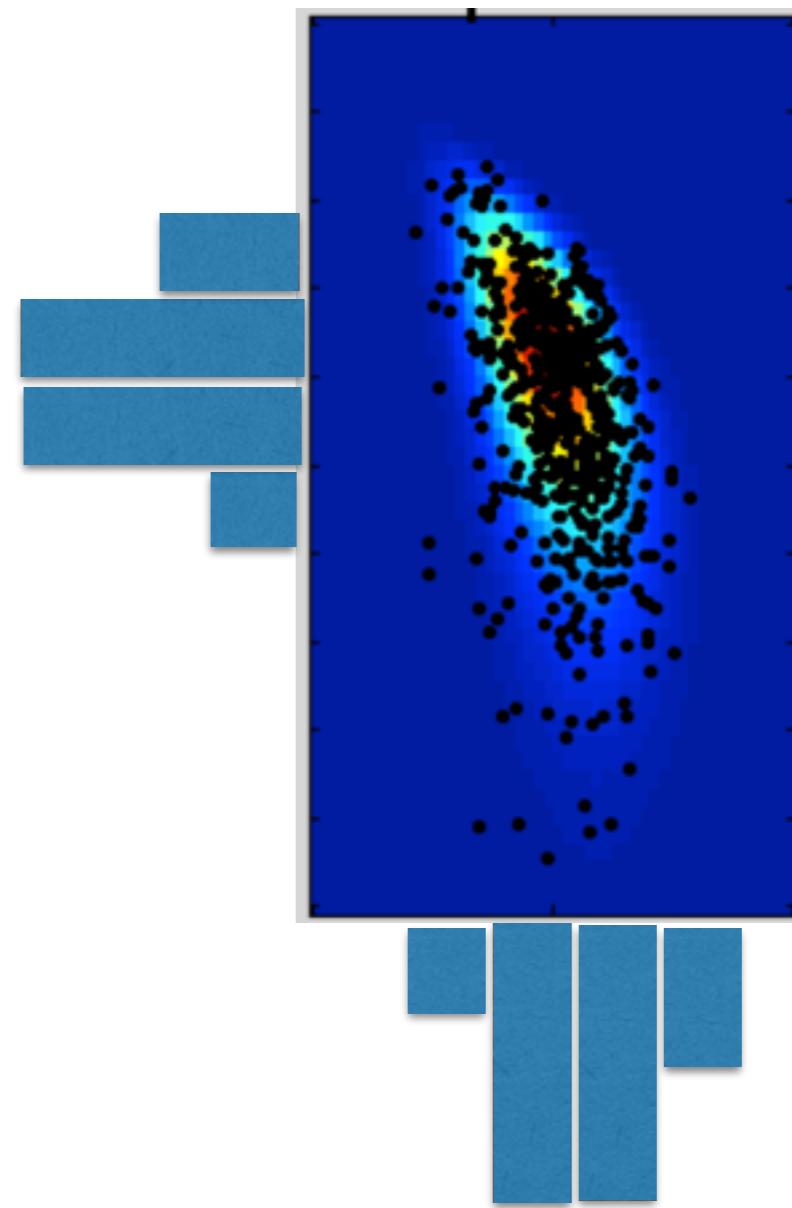
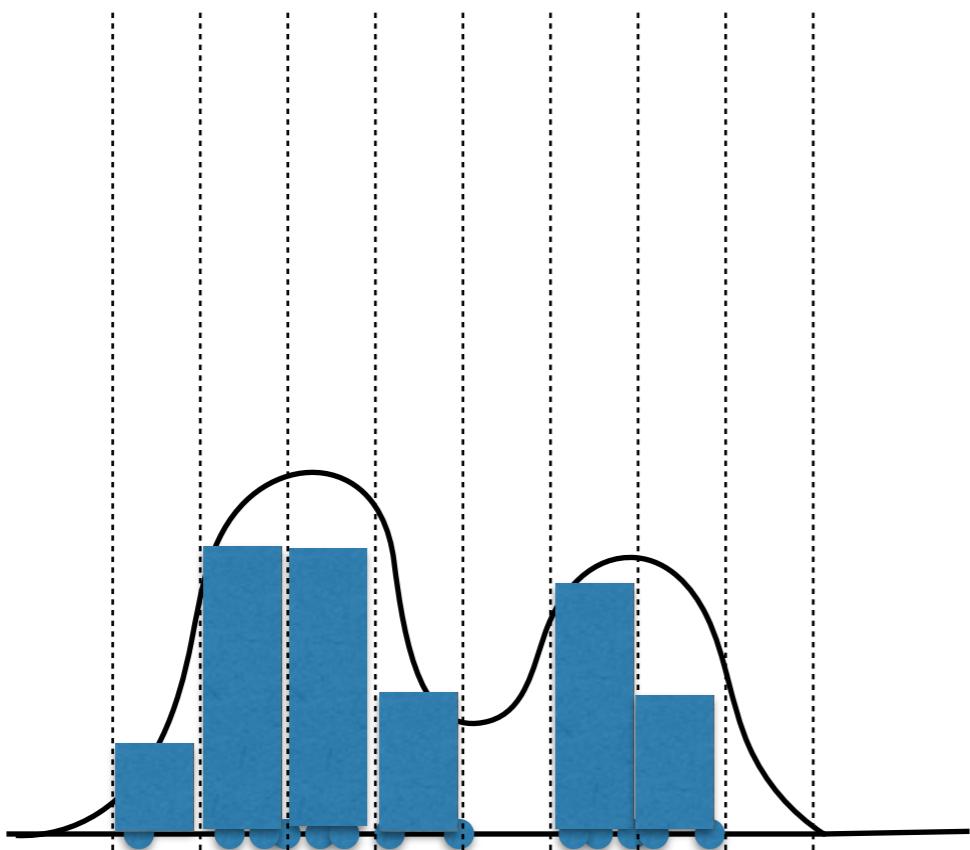
posterior

$$p(a|y) \propto N(a, m_p, s_p^2)$$

$$p(a|y) \propto |y-a|^{-n}$$



Sampling



Gibbs sampling

$$y = a^*x$$

noise model: $y = a^*x + N(0, s^2)$

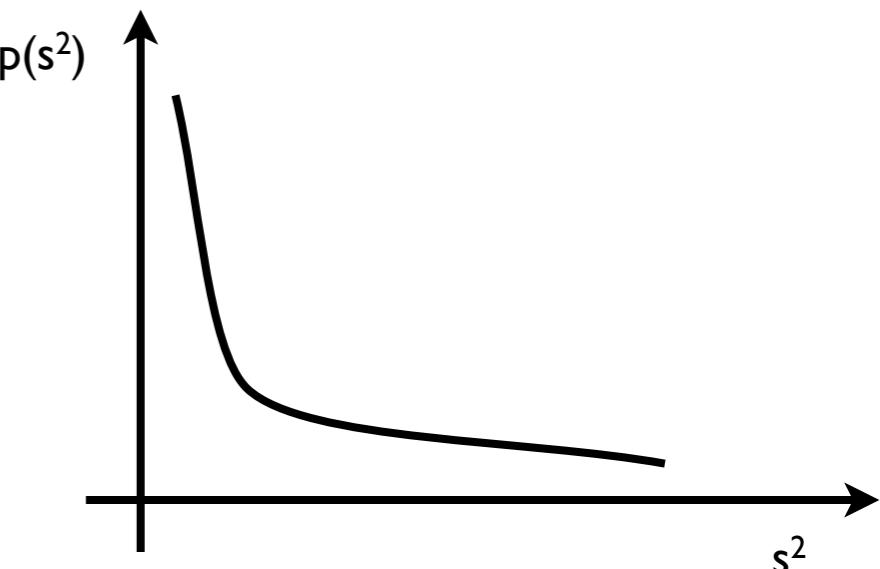
what if we want to learn about this one?

$$p(a, s^2 | y) \sim p(y|a, s^2) p(a) p(s^2)$$

likelihood: $p(y|a, s^2) = N(y|a^*x, s^2)$

$$p(a) = N(a|a_0, s_0^2)$$

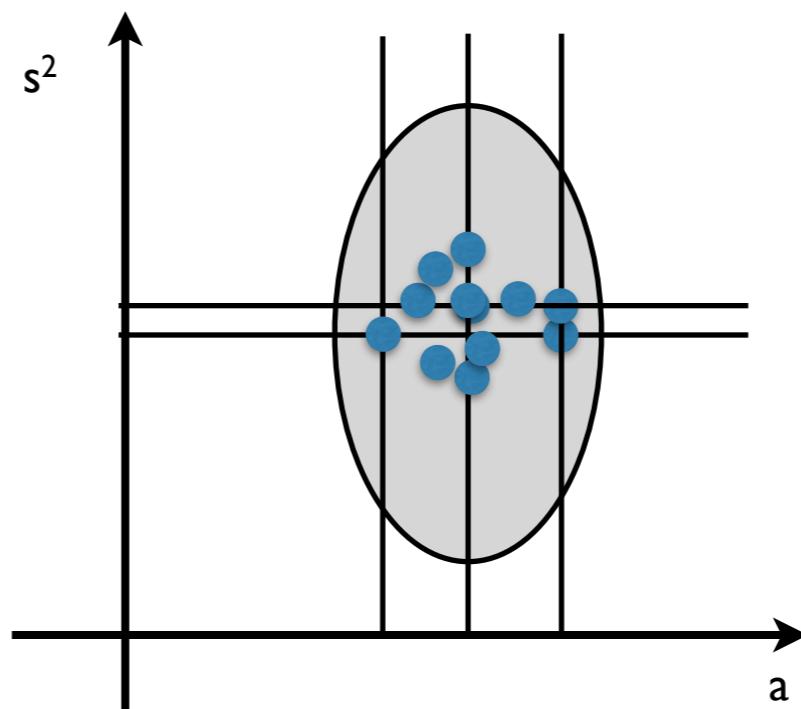
$$p(s^2) = 1/s^2$$



Gibbs sampling

iterate:

- sample from $p(a|y,s^2)$
- sample from $p(s^2|y,a)$



Gibbs sampling

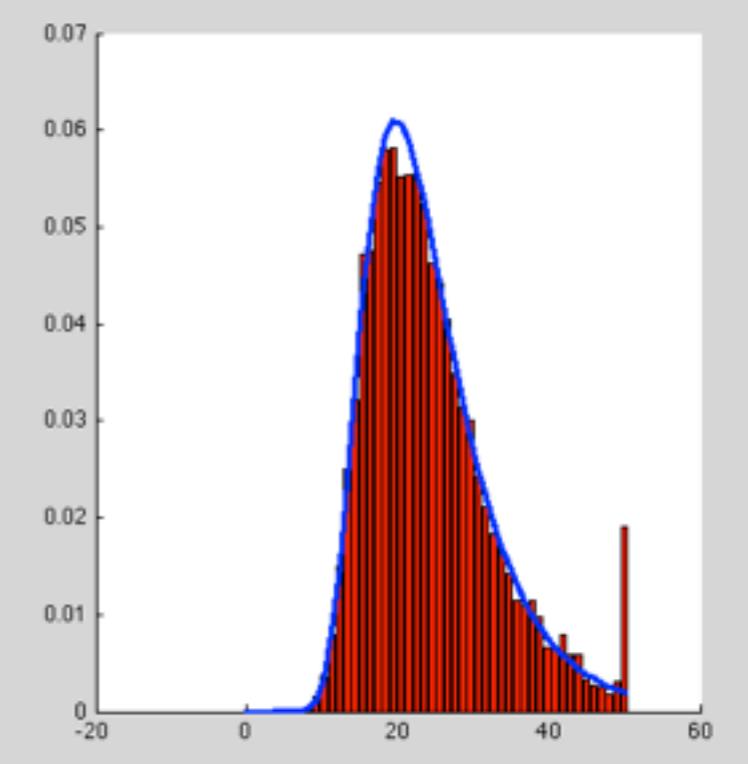
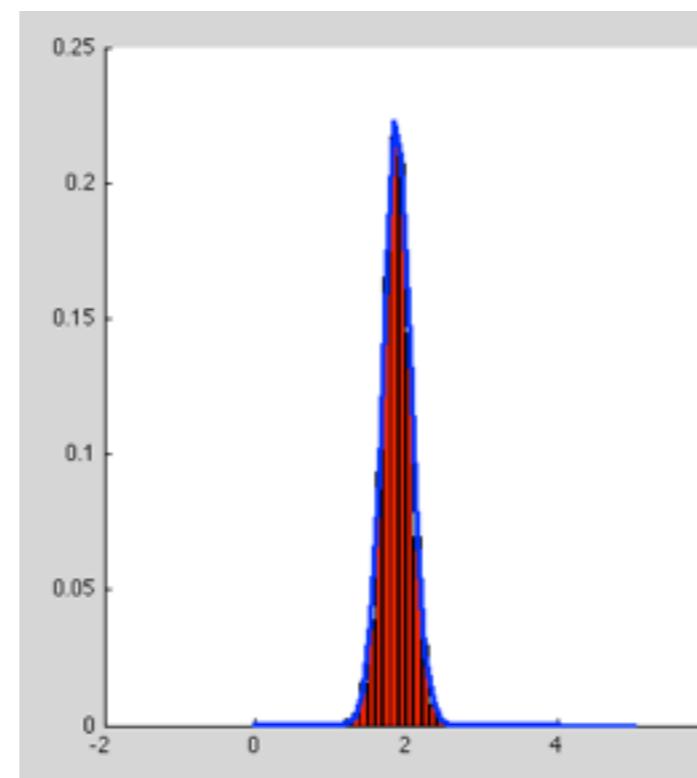
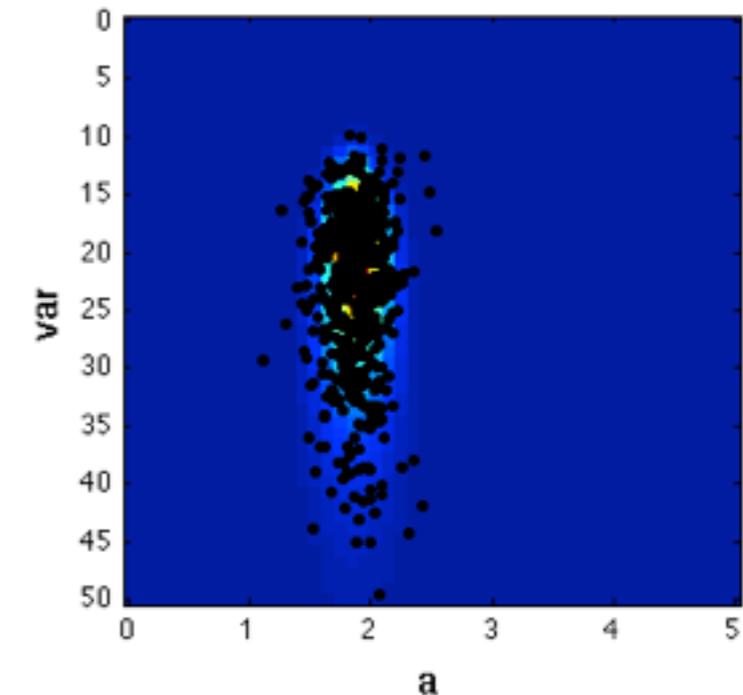
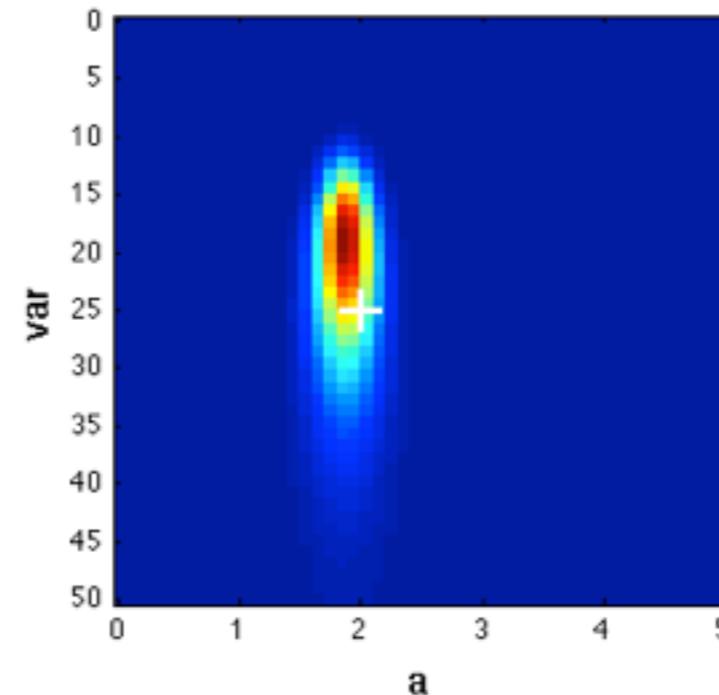
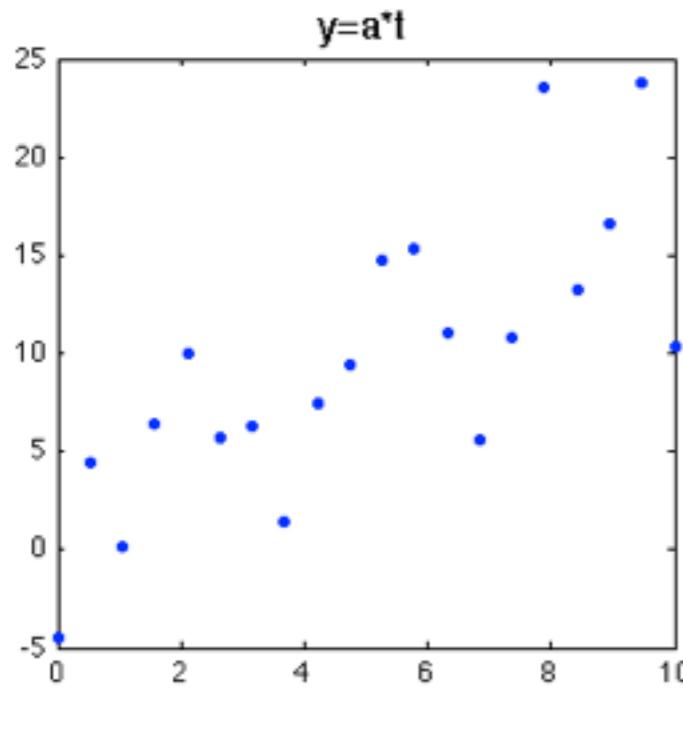
$p(a|y, s^2) = N(a_p, s_p)$ we can easily sample from this

$p(s^2|y, a) \propto p(y|a, s^2)p(s^2)$ product rule

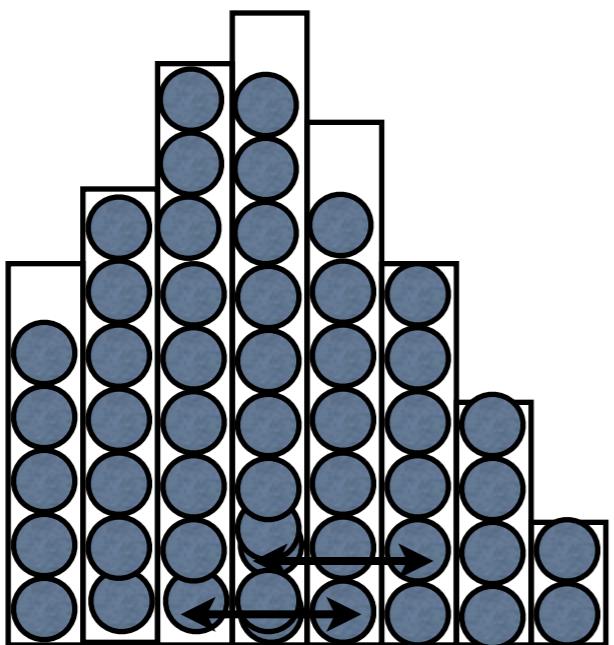
$$p(s^2|y, a) \propto (s^2)^{-n/2} \exp \frac{-\sum(y_i - ax_i)^2}{2s^2} \frac{1}{s^2}$$

$$\text{inv-G}(s^2|n/2, \sum(y_i - ax_i)^2))$$

Gibbs sampling



Metropolis Hastings



$a_1 \rightarrow a_2 \rightarrow a_3 \rightarrow a_4 \rightarrow a_4 \rightarrow a_5 \rightarrow a_6 \rightarrow a_7 \rightarrow a_8 \rightarrow a_9$
 q

$a_{n+1} \sim q(a_n)$ sample from q (typically Gaussian)

Metropolis Hastings

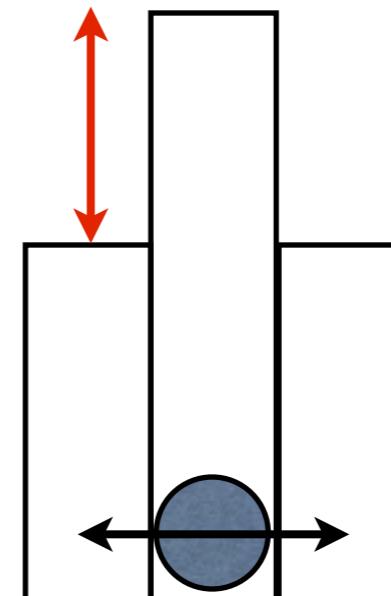
$a_{n+1} \sim q(a_n)$ sample from q (typically Gaussian)

if $p(a_{n+1}|y) > p(a_n|y)$
move!

if $p(a_{n+1}|y) < p(a_n|y)$

stay with probability proportional to ratio:

$$\frac{p(a_{n+1}|y)}{p(a_n|y)}$$



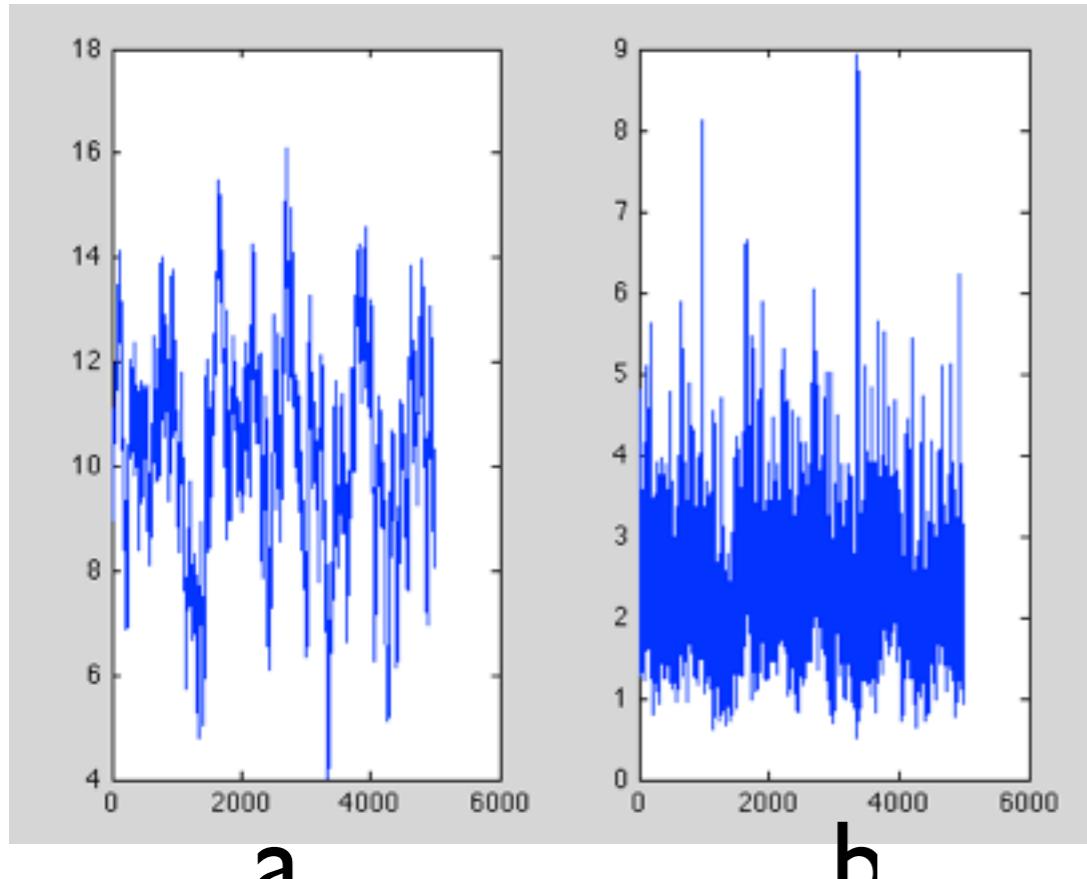
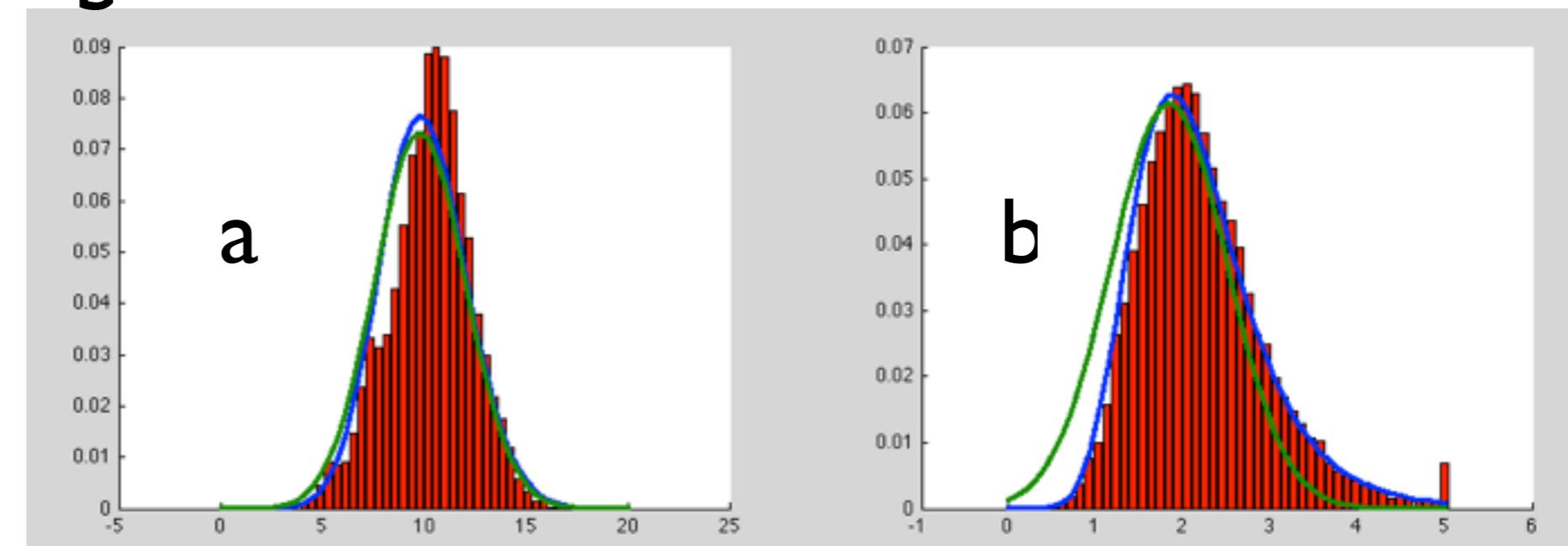
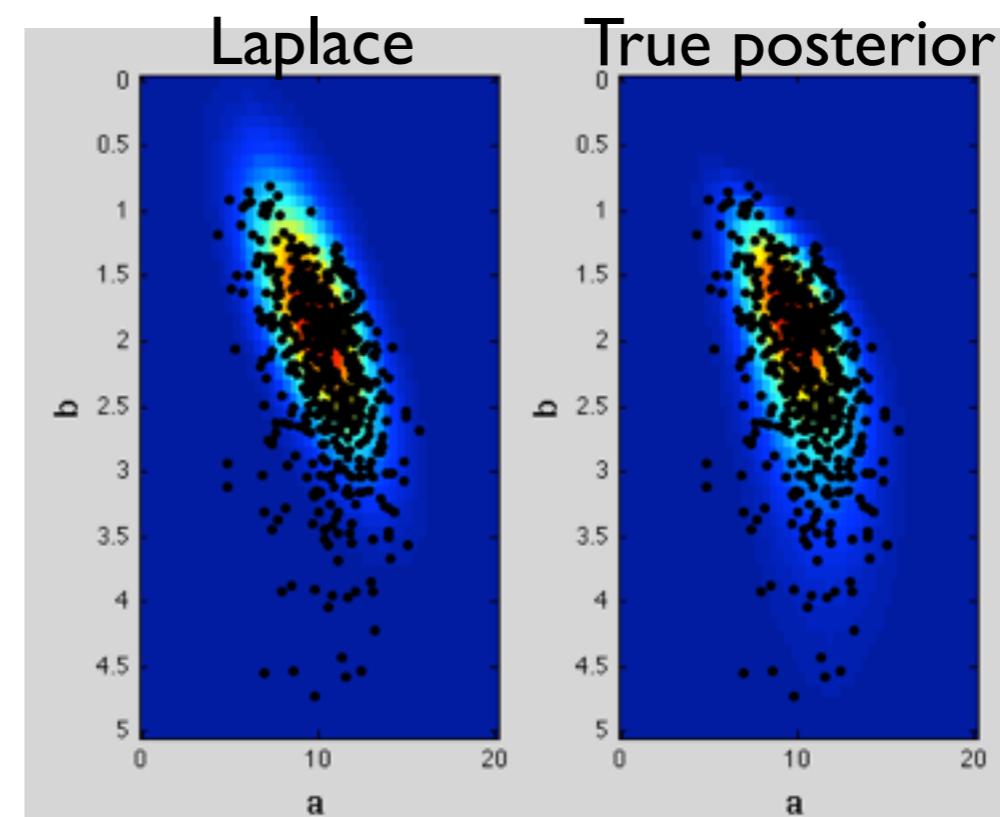
Final residence map
proportional to full posterior

Metropolis Hastings Code

```
function samples = mh(varargin)
% res = mh(y,x0,@genfunc,[LB,UB,params])
%
% Compulsory Parameters
%   y      = data (Nx1)
%   x0     = initial parameters (Px1)
%   @genfunc = generative model
%
% Output is nsamples*P where nsamples=njumps/sampleevery
%
% Example:
%   % define a forward model (here y=a*exp(-bx))
%   myfun=@(x,c)(exp(-x(1)*c)+x(2));
%   % generate some noisy data
%   true_x = [1;2];
%   c=linspace(1,10,100);
%   y=myfun(true_x,c) + .05*randn(1,100);
%   % estimate parameters
%   x0=[1;2]; % you can get x0 using nonlinear opt
%   samples=mh(y,x0,@(x)(myfun(x,c)));
%   figure,plot(samples)
%
% Other Parameters
%   LB = lower bounds on x (default=[-inf]*ones(size(x0)))
%   UB = upper bounds on x (default=[+inf]*ones(size(x0)))
%   params.burnin = #burnin iterations (default=1000)
%   params.njumps = #jumps (default=5000)
%   params.sampleevery = keep every n jumps (default=10)
%   params.update = update proposal every n jumps (default=20)
%
%
% S. Jbabdi 01/12
```

Metropolis Hastings

$$y = f(a, b) = ae^{-bt}$$

**a****b**

Metropolis Hastings

- extremely flexible
 - all sorts of priors possible
 - all sorts of models (linear/nonlinear)
 - high dimensional problems
- computationally intensive
- not easy to monitor convergence

Example

$$Y = M_0 * \exp(-TE/T_2) * (1 - \exp(-TR/T_1))$$

known

data, measurements
experimental parameters

unknown

model parameters

Example

$$Y = M_0 * \exp(-TE/T_2) * (1 - \exp(-TR/T_1))$$

$$Y_i = F_i(M_0, T_2, T_1)$$

Let's call these x

each 'i' is a different combination (TE,TR)

Example

$$Y = M_0 * \exp(-TE/T_2) * (1 - \exp(-TR/T_1))$$

$$Y_i = F_i(x)$$

“generative model”

Example

$$Y = M_0 * \exp(-TE/T_2) * (1 - \exp(-TR/T_1))$$

$$Y_i = F_i(x) + \text{noise}$$

noise trick:

$$p(x|Y) \propto |Y - F(x)|^{-n}$$

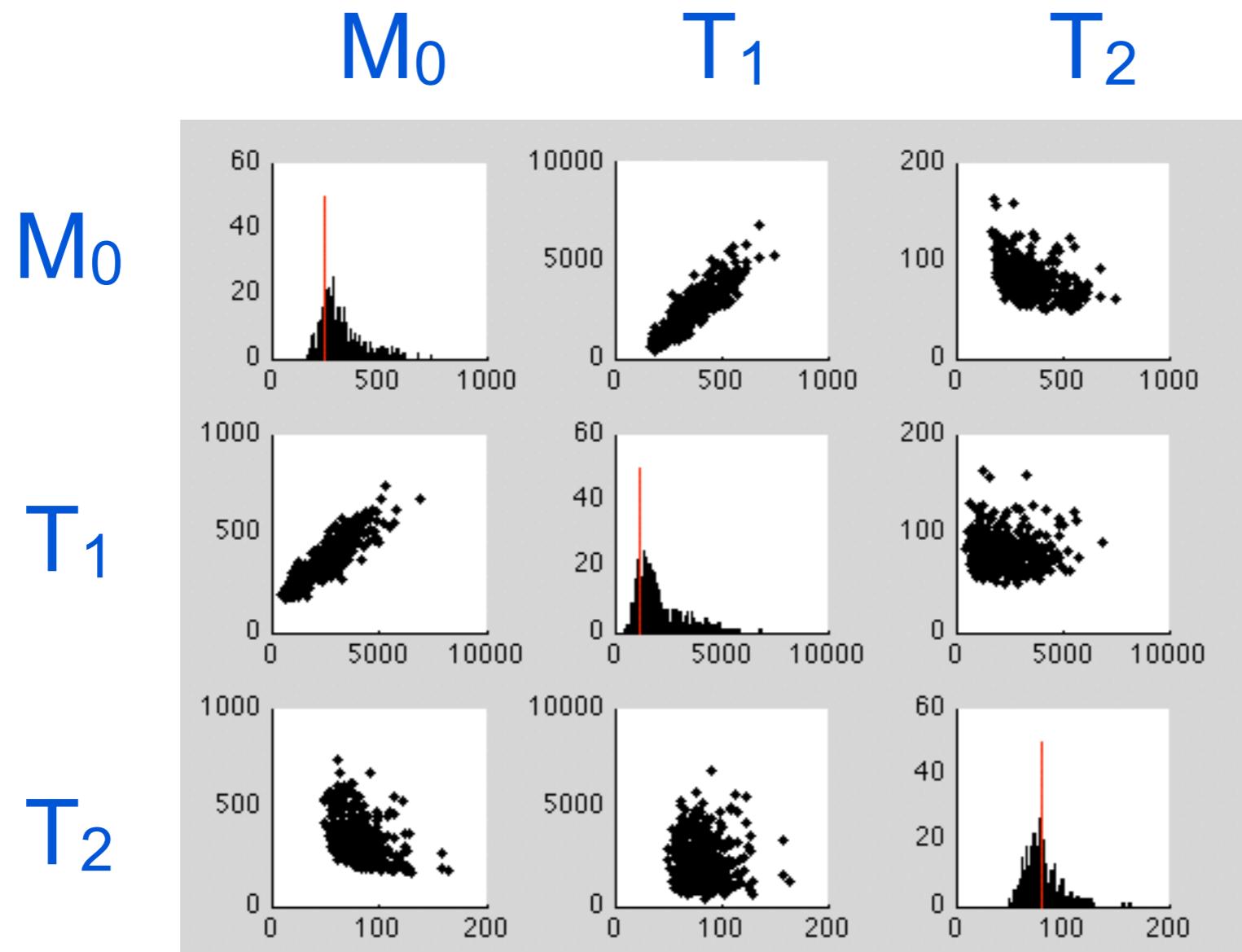
straight into MH!

TR = 1, 1.6, 2.3, 4 s

simple_mri_example.m

TE = 30, 47, 65, 82, 100 ms

$$Y = M_0 * \exp(-TE/T_2) * (1 - \exp(-TR/T_1))$$



SNR=20

**Lesson: Bayes can tell us about model structure (e.g.
correlation between parameters) and estimability
given data**

Bayesian Inference

- Exact inference

Analytic integrals (gaussian - or nasty)

Grid evalutation (1-5 parameters, max)

- Sampling

Metropolis Hastings (the most flexible, please use my tool!)

Gibbs sampling (need conditional posteriors to be nice)

- Analytic approximations

Laplace Approximation (local = curvature of posterior)

Variational Bayes (Kullback Leibler = approximate whole distribution)

ONBI - Bayesian modelling practical. 2014/15

Practical Overview

This practical requires Matlab. Go through the instructions and execute the listed commands in the Matlab command window (you can copy-paste). Raise your hand if you need help, but perhaps try first the "help" command in Matlab if you are unsure about Matlab syntax issues.

Contents:

- **Linear models**
Learn basics of fitting linear models
 - **Sampling**
Learn to use Gibbs sampling for inference
 - **Nonlinear models**
Learn to fit a nonlinear model using 3 different methods
-

Linear models

This practical will cover most of the examples that were discussed in the lecture, and you should be able to reproduce all the figures that were shown in the lecture.

We start with the simplest possible example of Bayesian inference. Imagine that we are measuring a certain quantity (e.g. temperature using a thermometer, actually let us say we ARE measuring temperature just for concreteness). We are going to assume that there exists a true unknown temperature, and that our measurement device is noisy. Therefore, we are going to make several measurements and we are going to use Bayesian inference.

We start by writing a generative model. Assuming we are measuring the temperature directly, the generative model can simply write:

```
y = a + noise;
```

In the above equation, y is the data, a is a parameter of the model (the true unknown temperature), and the noise is assumed to be additive.

We further need to make some distributional assumption on the noise. Most commonly, the noise is assumed Gaussian with mean zero and standard deviation s .