# SAMSI Astrostatistics Tutorial

## Models with Gaussian Uncertainties (lecture 1)

**Phil Gregory**

**University of British Columbia**

**2006**

# Bayesian Logical Data Analysis for the Physical Sciences

## A Comparative Approach with Mathematica Support

**P. C. Gregory**   (University of British Columbia, Vancouver)

**Hardback**         (ISBN-10: 052184150X | ISBN-13: 9780521841504)
**Also available in eBook format**

Increasingly, researchers in many branches of science are coming into contact with Bayesian statistics or Bayesian probability theory. By encompassing both inductive and deductive logic, Bayesian analysis can improve model parameter estimates by many orders of magnitude. It provides a simple and unified approach to all data analysis problems, allowing the experimenter to assign probabilities to competing hypotheses of interest, on the basis of the current state of knowledge. This book provides a clear exposition of the underlying concepts with large numbers of worked examples and problem sets. The book also discusses numerical techniques for implementing the Bayesian calculations, including an introduction to Markov Chain Monte-Carlo integration and linear and nonlinear least-squares analysis seen from a Bayesian perspective. In addition, background material is provided in appendices and supporting Mathematica notebooks are available, providing an easy learning route for upper-undergraduates, graduate students, or any serious researcher in physical sciences or engineering.

- Introduces statistical inference in the larger context of scientific methods, and includes many worked examples and problem sets.

- Presents Bayesian theory but also compares and contrasts with other existing ideas.

- Mathematica based support software is available.

Description
Table of contents
Excerpt
Index
Copyright
Frontmatter

**Resources and solutions
This title has free online support material available.**

View material...

Details
**128 line diagrams 4 half-tones 74 exercises 132 figures 55 worked examples**
Weight: 1.146 kg

# Bayesian Logical Data Analysis for the Physical Sciences

**Inference with Gaussian uncertainties**

# Resources and solutions

**Book Preface** (11 Kb)

**Mathematica Tutorial** (1415 Kb)

**Errata** (206 Kb)

**Revisions** (268 Kb)

Book website:   www.cambridge.org/052184150X

# Outline

# How to proceed in a Bayesian analysis?

**Write down Bayes' theorem, identify the terms and solve.**

**Prior probability**      **Likelihood**

$$p(H_i \,/D,I) \;=\; \frac{p(H_i \mid I) \times p(D \mid H_i, I)}{p(D \mid I)}$$

**Posterior probability that $H_i$ is true, given the new data D and prior information I**

**Normalizing constant**

**Every item to the right of the vertical bar | is assumed to be true**

**The likelihood $p(D \mid H_i, I)$, also written as $\mathcal{L}(H_i)$, stands for the probability that we would have gotten the data D that we did, if $H_i$ is true.**

# Simple Spectral Line Problem

## Background (prior) information:

**Two competing grand unification theories have been proposed, each championed by a Nobel prize winner in physics. We want to compute the relative probability of the truth of each theory based on our prior information and some new data.**

**Theory 1 is unique in that it predicts the existence of a new short-lived baryon which is expected to form a short-lived atom and give rise to a spectral line at an accurately calculable radio wavelength.**

**Unfortunately, it is not feasible to detect the line in the laboratory. The only possibility of obtaining a sufficient column density of the short-lived atom is in interstellar space.**

**Prior estimates of the line strength expected from the Orion nebula according to theory 1 range from 0.1 to 100 mK.**

# Simple Spectral Line Problem

**The predicted line shape has the form**

$$T \exp \left\{ \frac{-(\nu_i - \nu_o)^2}{8} \right\} \quad \text{(abbreviated by } Tf_i\text{)},$$

**where the signal strength is measured in temperature units of mK and $T$ is the amplitude of the line. The frequency, $\nu_i$ , is in units of the spectrometer channel number and the line center frequency $\nu_0$ = 37.**



Line profile $f_i$

# Data

**To test this prediction, a new spectrometer was mounted on the James Clerk Maxwell telescope on Mauna Kea and the spectrum shown below was obtained. The spectrometer has 64 frequency channels with neighboring.**



**All channels have Gaussian noise characterized by $\sigma$ = 1 mK. The noise in separate channels is independent. The line center frequency $\nu_0$ = 37.**

# Questions of interest

**Based on our current state of information, which includes just the above prior information and the measured spectrum,**

**1) what do we conclude about the relative probabilities of the two competing theories**

*and*

**2) what is the posterior PDF for the line strength?**

**Hypothesis space of interest for model selection part:**

$$M_1 \equiv \text{``Theory 1 correct, line exists''}$$

$$M_2 \equiv \text{``Theory 2 correct, no line predicted''}$$

# Model selection

**To answer the model selection question, we compute the odds ratio (abbreviated simply by the *odds*) of model $M_1$ to model $M_2$ .**

**Expand numerator and denominator with Bayes' theorem**

$$O_{12} = \frac{p(M_1 \mid D, I)}{p(M_2 \mid D, I)} = \frac{\frac{p(M_1 \mid I)\, p(D \mid M_1, I)}{p(D \mid I)}}{\frac{p(M_2 \mid I)\, p(D \mid M_2, I)}{p(D \mid I)}} = \frac{p(M_1 \mid I)}{p(M_2 \mid I)} \frac{p(D \mid M_1, I)}{p(D \mid M_2, I)}$$

**posterior probability ratio**

**prior probability ratio**

**Bayes factor**

**$p(D|M_1, I)$, the called the global likelihood of $M_1$ .**

$$p(D \mid M_1, I) = \int_T p(D, T \mid M_1, I)\, dT$$

**Expanded with product rule**

$$= \int_T p(T \mid M_1, I)\, p(D \mid M_1, T, I)\, dT$$

**The global likelihood of a model is equal to the weighted average likelihood for its parameters.**

# Choice of prior $p(T|M_1, I)$

## Investigate two common choices

1. Uniform prior
$$p(T|M_1, I) = \frac{1}{\Delta T}$$

where $\Delta T = T_{\max} - T_{\min}$

**There is a problem with this prior if the range of *T* is large. In the current example $T_{\min}$ = 0.1 and $T_{\max}$ = 100. Compare the probability that *T* lies in the upper decade of the prior range (10 to 100 mK) to the lowest decade (0.1 to 1 mK).**

$$\frac{\int_{10}^{100} p(T|M_1, I)dT}{\int_{0.1}^{1} p(T|M_1, I)dT} = 100$$

**Usually, expressing great uncertainty in some quantity corresponds more closely to a statement of scale invariance or equal probability per decade. The Jeffreys prior, discussed next, has this scale invariant property.**

# Choice of prior $p(T|M_1, I)$

2. **Jeffreys prior**
   **(Scale invariant)**

$$p(T|M_1, I) = \frac{1}{T \, \ln\left(\frac{T_{max}}{T_{min}}\right)}$$

$$\int_{0.1}^{1} p(T|M_1, I)dT = \int_{10}^{100} p(T|M_1, I)dT$$

**What if the lower bound on T includes zero? Another alternative Is a modified Jeffreys prior of the form.**

$$p(T \mid M_1, I) = \frac{1}{T + T_0} \frac{1}{\ln\left(\frac{T_0 + T_{max}}{T_0}\right)}$$

**This prior behaves like a uniform prior for $T < T_0$ and a Jeffreys prior for $T > T_0$. Typically set $T_0$ = noise level.**

## Calculation of $p(D|M_1, T, I)$

**Let $d_i$ represent the measured data value for the $i^{th}$ channel of the spectrometer. According to model $M_1$,**

$$d_i = T f_i + e_i \quad \text{and} \quad f_i = \exp\left(\frac{-(\nu_i - \nu_0)^2}{8}\right),$$

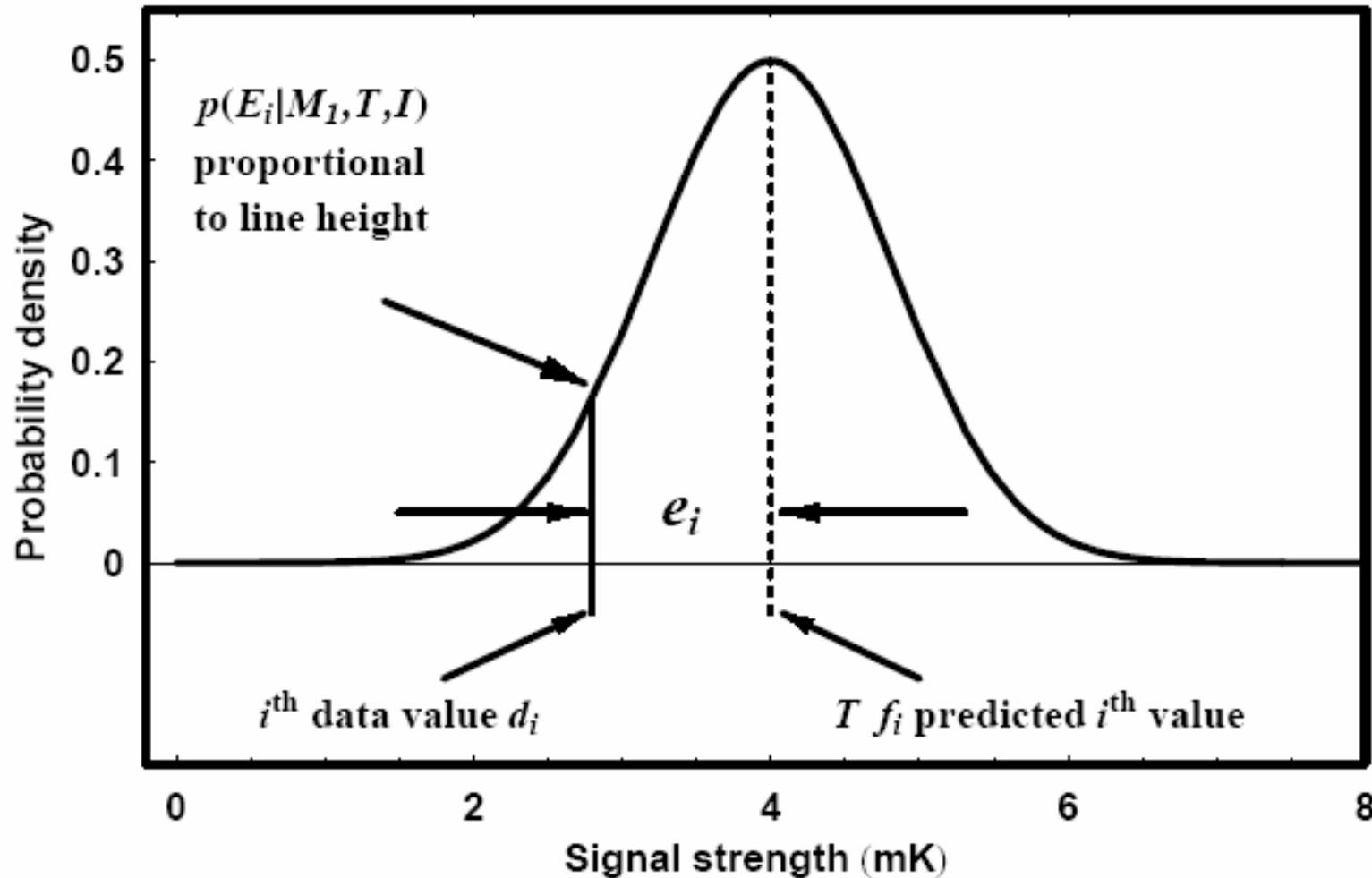**and $e_i$ represents the error component in the measurement. Our prior information indicates that $e_i$ has a Gaussian distribution with a $\sigma = 1\ mK$.**

**Assuming $M_1$ is true, then if it were not for the error $e_i$, $d_i$ would equal $T f_i$.**

**Let $E_i \equiv$ "a proposition asserting that the $i^{th}$ error value is in the range $e_i$ to $e_i + de_i$."    If all the $E_i$ are independent  then**

$$
\begin{aligned}
p(D|M_1, T, I) &= p(D_1, D_2, ..., D_N | M_1, T, I) \\
&= p(E_1, E_2, ..., E_N | M_1, T, I) \\
&= p(E_1 | M_1, T, I) p(E_2 | M_1, T, I) ... p(E_N | M_1, T, I) \\
&= \prod_{i=1}^{N} p(E_i | M_1, T, I)
\end{aligned}
$$

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# Calculation of $p(D|M_1, T, I)$

**Probability of getting a data value $d_i$ a distance $e_i$ away from the predicted value is proportional to the height of the Gaussian error curve at that location.**

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# Calculation of $p(D|M_1, T, I)$

**From the prior information, we can write**

$$
\begin{aligned}
p(E_i|M_1, T, I) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{e_i^2}{2\sigma^2}\right\} \\
&= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(d_i - Tf_i)^2}{2\sigma^2}\right\}
\end{aligned}
$$

**Our final likelihood is given by**

$$
\begin{aligned}
p(D|M_1, T, I) &= \prod_{i=1}^{N} \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(d_i - Tf_i)^2}{2\sigma^2}\right) \\
&= (2\pi)^{\frac{-N}{2}} \sigma^{-N} \exp\left\{-\frac{\sum_i (d_i - Tf_i)^2}{2\sigma^2}\right\}
\end{aligned}
$$

# Calculation of $p(D|M_2, I)$

**Model $M_2$ assumes the spectrum is consistent with noise and has no free parameters so we can write**

$$d_i = 0 + e_i$$

$$p(D|M_2, I) = (2\pi)^{-\frac{N}{2}} \sigma^{-N} e^{-\frac{\sum d_i^2}{2\sigma^2}}$$

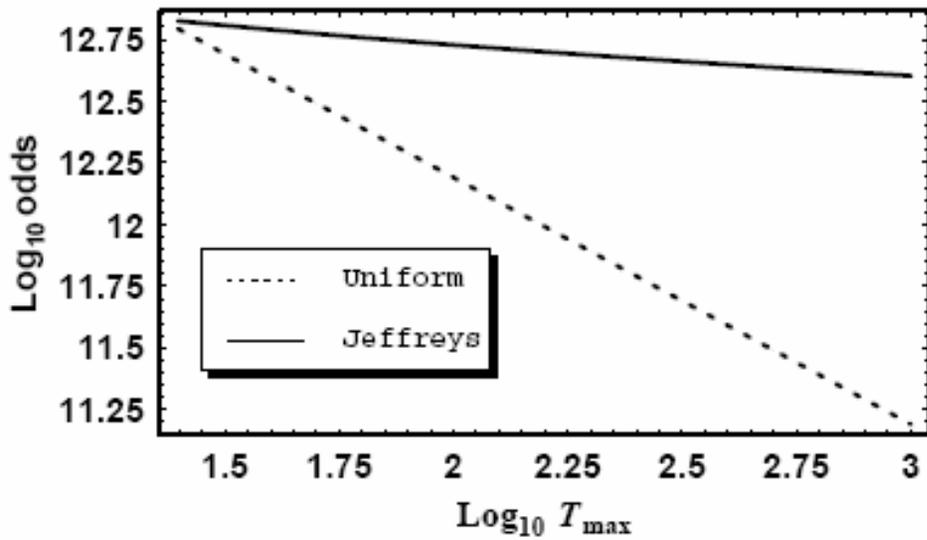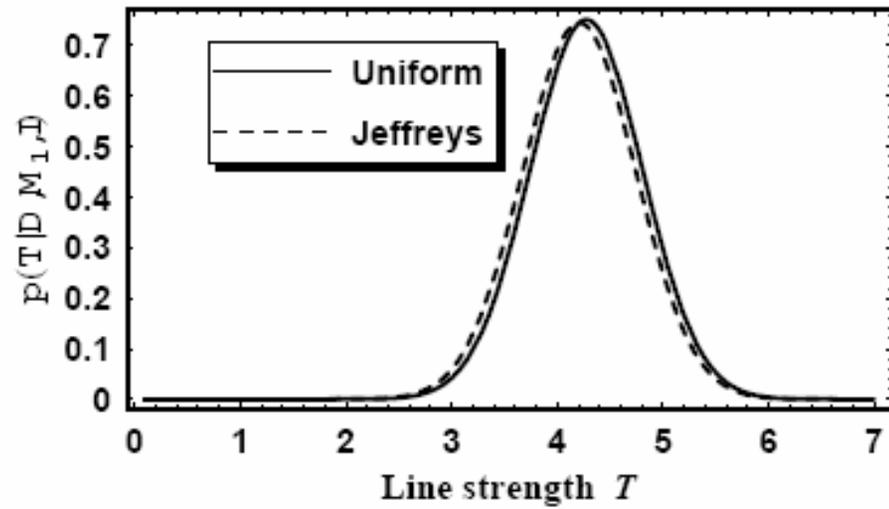## Model selection results

**Bayes factor, uniform prior = $1.6 \times 10^{12}$**

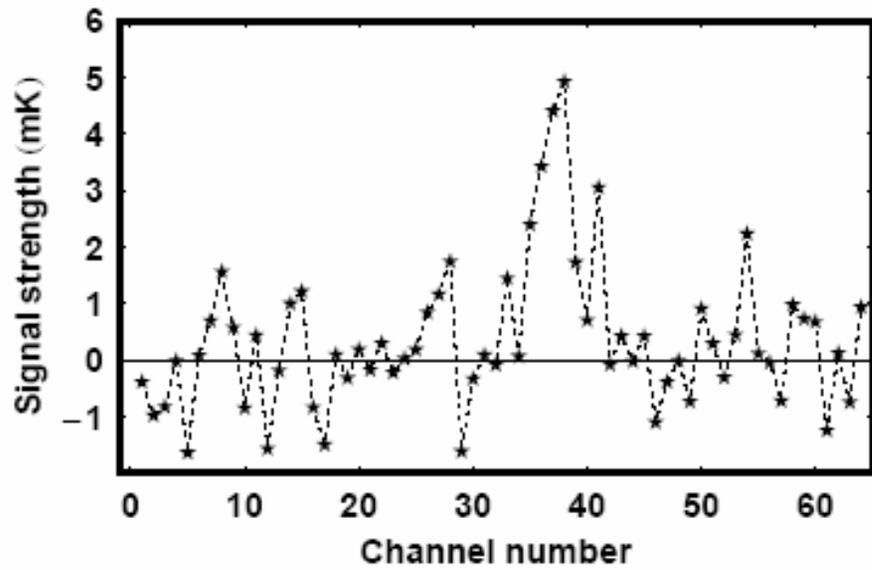**Bayes factor, Jeffreys prior = $5.3 \times 10^{12}$**

**The factor of $10^{12}$ is so large that we are not terribly interested in whether the factor in front is 1.6 or 5.3. Thus the choice of prior is of little consequence when the evidence provided by the data for the existence is as strong as this.**

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005
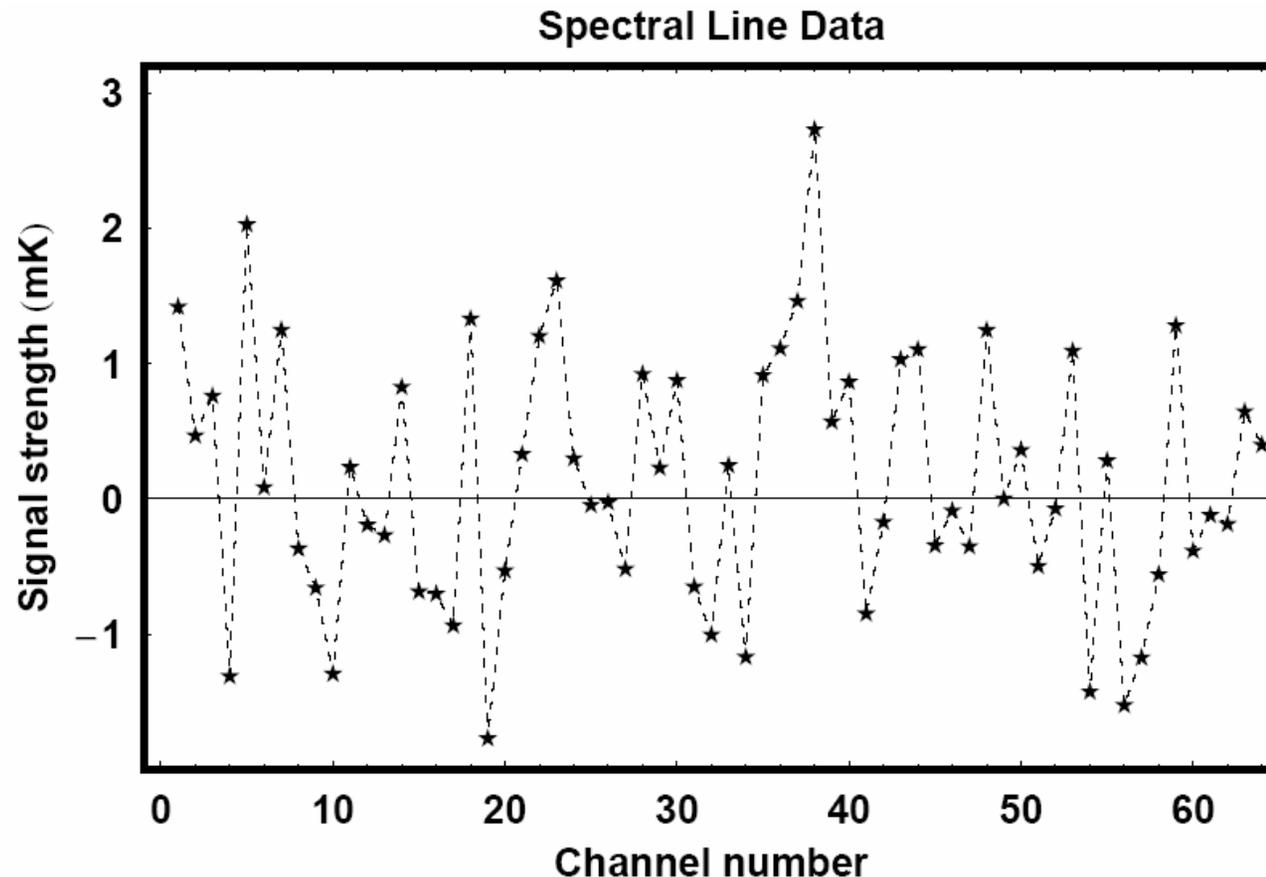
# Parameter Estimation Problem

**Now that we have solved the model selection problem leading to a significant preference for $M_1$, we would now like to compute $p(T|D,M_1, I)$, the posterior PDF for the signal strength.**

Again, start with Bayes' theorem

$$
\begin{aligned}
p(T|D, M_1, I) &= \frac{p(T|M_1, I)p(D|M_1, T, I)}{p(D|M_1, I)} \\
&\propto p(T|M_1, I)p(D|M_1, T, I)
\end{aligned}
$$

# How do our conclusions change when evidence for the line in the data is weaker?


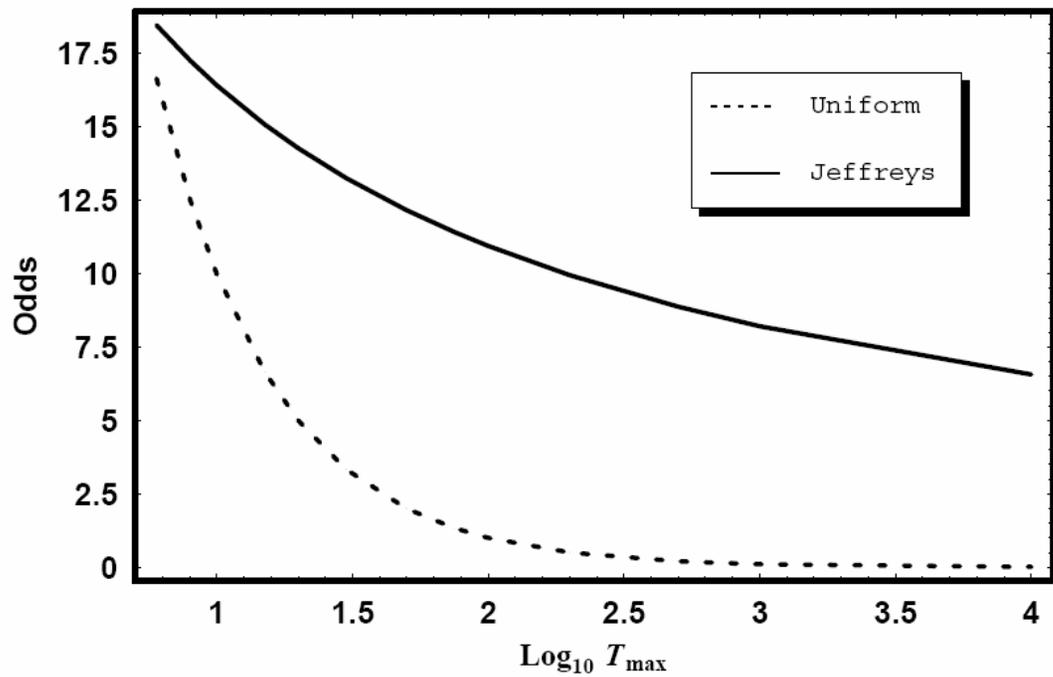
**Spectral Line Data**

**All channels have IID Gaussian noise characterized by $\sigma$ = 1 mK.**
**The predicted line center frequency $\nu_0$ = 37.**

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# Model selection results

**Bayes factor, uniform prior = 1.0**

**Bayes factor, Jeffreys prior = 11.**

**As expected, when the evidence provided by the data is much weaker, our conclusions can be strongly influenced by the choice of prior and it is a good idea to examine the sensitivity of the results by employing several different priors.**

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# What if we were uncertain about the line center frequency?

Suppose our prior information only restricted the center frequency to the first 44 channels of the spectrometer.

In this case $\nu_0$ becomes a nuisance parameter that we can marginalize.

$$p(D \mid M_1, I) = \int_{\nu_0} \int_T p(D, T, \nu_0 \mid M_1, I)\, dT\, d\nu_0$$

Assumes independent priors

$$= \int_{\nu_0} \int_T p(T \mid M_1, I)\, p(\nu_0 \mid M_1, I)\, p(D \mid M_1, T, \nu_0, I)\, dT\, d\nu_0$$

**New Bayes factor = 1.0 , assuming a uniform prior for $\nu_0$ and a Jeffreys prior for T.**

Built into any Bayesian model selection calculation is an automatic and quantitative Occam's razor that penalizes more complicated models. One can identify an Occam factor for each parameter that is marginalized. The size of any individual Occam factor depends on our prior ignorance in the particular parameter.

# Marginal PDF for line center frequency



**Marginal posterior PDF for the line frequency, where the line frequency is expressed as a spectrometer channel number.**

# Marginal probability density function for line strength

**Line frequency (uniform prior) & line strength (Jeffreys)**



Legend:
- $\nu_0$ unknown
- $\nu_0$ = channel 37

Y-axis: PDF

X-axis: Line Strength (mK)

# Joint probability density function

Contours$\rightarrow${90, 50, 20, 10, 5, 2, 1, 0.5}% of peak



Line strength (mK) vs Frequency channels (1 to 44)

# Generalizations

$$d_i = T \times \exp\left(-\frac{(\nu_i - \nu_0)^2}{8}\right) + e_i$$

<span style="color:red">**current model**</span>

**More generally we can write**

$$d_i = f_i + e_i$$

**where** $\quad f_i = \sum_{\alpha=1}^{m} A_\alpha \, g_\alpha(x_i)$

**specifies a linear model with m basis functions** $g_\alpha(x_i)$

**or** $\quad f_i = \sum_{\alpha=1}^{m} A_\alpha \, g_\alpha(x_i \mid \theta)$

**specifies a model with m basis functions with an additional set of nonlinear parameters represented by $\theta$.**

## Generalizations

**Examples of linear models**

$$f_i = A_1 + A_2\,x + A_3\,x^2 + A_3\,x^3 + \dots = \sum_{\alpha=1}^{m} A_\alpha\,g_\alpha\,(x_i)$$

$$f_i = A_1 \times \exp\left(-\frac{(\nu_i - C_1)^2}{2\,\sigma_1^2}\right) \quad \text{where } C_1 \text{ and } \sigma_1 \text{ are known.}$$

**Examples of nonlinear models**

$$f_i = A_1\,\cos\omega t + A_2\,\sin\omega t$$

where $A_1, A_2$ and $\omega$ are unknowns.

$$f_i = A_1 \times \exp\left(-\frac{(\nu_i - C_1)^2}{2\,\sigma_1^2}\right) + A_2 \times \exp\left(-\frac{(\nu_i - C_2)^2}{2\,\sigma_2^2}\right) + \dots$$

where the A's, C's and $\sigma$'s are unknowns.

**Data** **Model** **Prior**
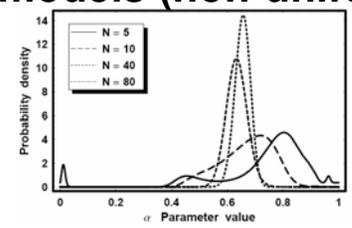**D** **M** **I**

**Posterior**

**Linear models** (uniform priors)

**Nonlinear models**

**+ linear models (non-uniform priors)**



**Posterior has a single peak
(multi-dimensional Gaussian)**

**Parameters given
by the normal equations
of linear least-squares**

**No integration required
very fast using
linear algebra**

**Posterior may have multiple peaks**

| Brute force integration | Asymptotic approx.'s | Moderate dimensions | High dimensions |
|---|---|---|---|
| **For some nuisance parameters analytic integration sometimes possible** | peak finding algorithms (e.g. Levenberg-Marquardt) | quadrature | MCMC |
| | Laplace approx.'s | randomized quadrature | |
| | | adaptive quadrature | |

**(chapter 10)**          **(chapter 11)**          **(chapter 12)**

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

**The Bayesian posterior density function for the nonlinear model parameter α for 4 simulated data sets of different size ranging from *N* = 5 to *N* = 80. The *N* = 5 case has the broadest distribution and exhibits 4 maxima.**

# Mean: Known noise $\sigma$

The problem is to solve for $p(\mu|D, I)$. The first step is to write down Bayes' theorem:

$$p(\mu|D, I) = \frac{p(\mu|I)\, p(D|\mu, I)}{p(D|I)},$$

$$
\begin{aligned}
p(\mu|I) &= K \text{ (constant)} & \mu_{\mathrm{L}} \leq \mu \leq \mu_{\mathrm{H}} \\
&= 0 & \text{otherwise.}
\end{aligned}
$$

The likelihood is given by

$$p(D|\mu, I) = \prod_{i=1}^{N} \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(d_i - \mu)^2}{2\sigma^2}\right\}$$

**Answer**

$$p(\mu|D, I) = \frac{\exp\left\{-\frac{(\mu - \bar{d})^2}{2\sigma^2/N}\right\}}{\int_{\mu_{\mathrm{L}}}^{\mu_{\mathrm{H}}} \exp\left\{-\frac{(\mu - \bar{d})^2}{2\sigma^2/N}\right\} d\mu}$$

# Mean: Known noise, unequal $\sigma$

$$p(D|\mu, I) = \prod_{i=1}^{N} \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(d_i - \mu)^2}{2\sigma_i^2}}$$

$$= \left[\prod_{i=1}^{N} \sigma_i^{-1}\right] (2\pi)^{-\frac{N}{2}} \exp\left\{-\sum_{i=1}^{N} \frac{(d_i - \mu)^2}{2\sigma_i^2}\right\}$$

$$= \left[\prod_{i=1}^{N} \sigma_i^{-1}\right] (2\pi)^{-\frac{N}{2}} \exp\left\{-\sum_{i=1}^{N} \frac{w_i(d_i - \mu)^2}{2}\right\}$$

where $w_i = 1/\sigma_i^2$ is called the weight of data value $d_i$.

**Answer**

$$p(\mu|D, I) = \frac{\exp\left\{-\frac{(\mu - \overline{d_w})^2}{2\sigma_w^2}\right\}}{\int_{\mu_L}^{\mu_H} \exp\left\{-\frac{(\mu - \overline{d_w})^2}{2\sigma_w^2}\right\} d\mu}.$$

Since the denominator evaluates to a constant, the posterior, within the range $\mu_L$ to $\mu_H$, is simply a Gaussian with variance $\sigma_w^2 = 1/(\sum w_i)$. The most probable value of $\mu$ is the weighted mean $\overline{d_w} = \sum w_i d_i / (\sum w_i)$.

**Mean: Unknown noise $\sigma$**



It is obvious from the scatter in the measurements compared to the error bars that there is some additional source of uncertainty or the signal strength is variable. For example, additional fluctuations might arise from propagation effects in the interstellar medium between the source and observer.

In the absence of prior information about the distribution of the additional scatter, both the **Central Limit Theorem** and the **Maximum Entropy Principle** lead us to adopt a Gaussian distribution because it is the most conservative choice.

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# Mean: Unknown noise $\sigma$

Now we have two unknowns in our model, $\mu$ and $\sigma$. The joint posterior probability $p(\mu, \sigma | D, I)$ is given by Bayes' theorem:

$$p(\mu, \sigma | D, I) = \frac{p(\mu, \sigma | I) p(D | \mu, \sigma, I)}{p(D | I)}.$$

We are interested in $p(\mu | D, I)$ regardless of what the true value of $\sigma$ is. In this problem, $\sigma$ is a nuisance parameter so we marginalize over $\sigma$:

$$p(\mu | D, I) = \int p(\mu, \sigma | D, I) d\sigma.$$

From the product rule: $p(\mu, \sigma | I) = p(\mu | I) \, p(\sigma | \mu, I)$.
Assuming the prior for $\sigma$ is independent of the prior for $\mu$, then

$$p(\mu, \sigma | I) = p(\mu | I) \, p(\sigma | I).$$

# Mean: Unknown noise $\sigma$

We will assume a Jeffreys prior for the scale parameter $\sigma$:

$$p(\sigma|I) = \begin{cases} \frac{K}{\sigma} & \sigma_{\mathrm{L}} \leq \sigma \leq \sigma_{\mathrm{H}} \\ 0 & \text{otherwise.} \end{cases}$$

The constant $K$ is determined from the condition

$$\int_{\sigma_{\mathrm{L}}}^{\sigma_{\mathrm{H}}} p(\sigma|I)d\sigma = 1 \quad \Rightarrow \quad K = \frac{1}{\ln(\frac{\sigma_{\mathrm{H}}}{\sigma_{\mathrm{L}}})}$$

$$p(\sigma|I) = \frac{1}{\sigma \ln \frac{\sigma_{\mathrm{H}}}{\sigma_{\mathrm{L}}}}. \tag{9.26}$$

**After some maths and a change of variables we can write**

$$p(\mu|D,I) = \frac{Q^{-(\frac{N}{2})} \int_{\tau_{\mathrm{L}}}^{\tau_{\mathrm{H}}} \tau^{\frac{N}{2}-1} e^{-\tau} d\tau}{\int_{\mu_{\mathrm{L}}}^{\mu_{\mathrm{H}}} d\mu \, Q^{-(\frac{N}{2})} \int_{\tau_{\mathrm{L}}}^{\tau_{\mathrm{H}}} \tau^{\frac{N}{2}-1} e^{-\tau} d\tau}$$

$$\approx \frac{Q^{-(\frac{N}{2})}}{\int_{\mu_{\mathrm{L}}}^{\mu_{\mathrm{H}}} d\mu \, Q^{-(\frac{N}{2})}}, \tag{9.34}$$

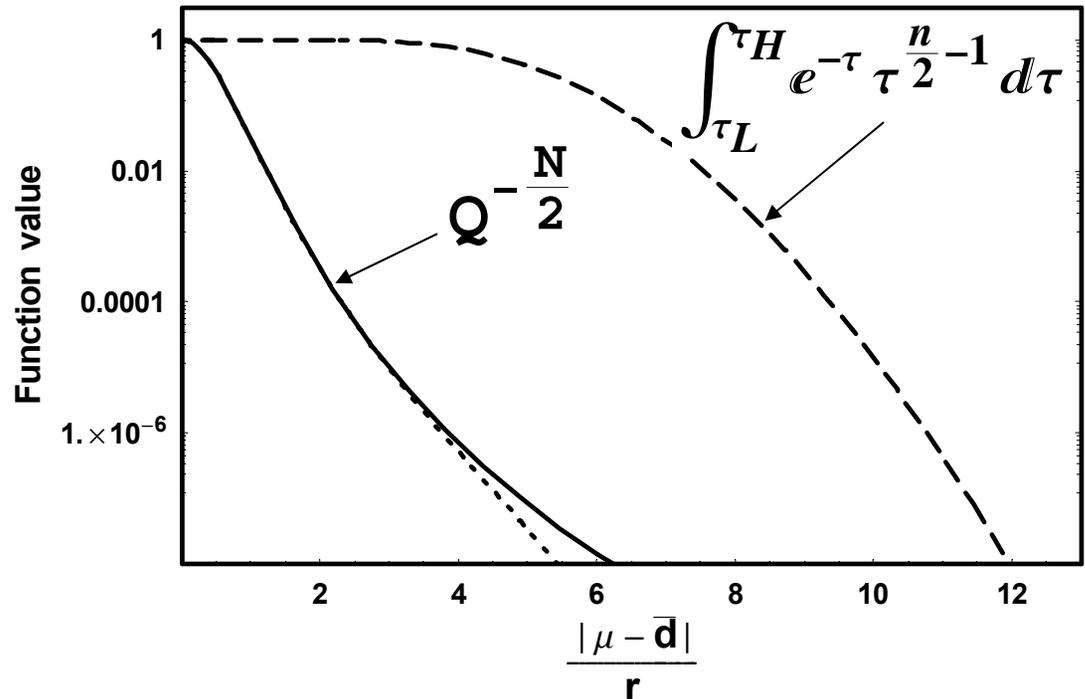**where** $Q = N(\mu - \overline{d})^2 + Nr^2$, $\quad Nr^2 = \sum(d_i - \overline{d})^2$, **and** $\tau = \frac{Q}{2\sigma^2}$

# Justification for dropping terms involving $\int_{\tau_L}^{\tau_H} \tau^{\frac{N}{2}-1} e^{-\tau} d\tau$

$$p(\mu|D,I) \;=\; \frac{Q^{-(\frac{N}{2})} \int_{\tau_L}^{\tau_H} \tau^{\frac{N}{2}-1} e^{-\tau} d\tau}{\int_{\mu_L}^{\mu_H} d\mu \, Q^{-(\frac{N}{2})} \int_{\tau_L}^{\tau_H} \tau^{\frac{N}{2}-1} e^{-\tau} d\tau}$$

$$\approx \; \frac{Q^{-(\frac{N}{2})}}{\int_{\mu_L}^{\mu_H} d\mu \, Q^{-(\frac{N}{2})}},$$

N = 10,  $\sigma_H$ = 5 r,  $\sigma_L$ = 0.5 r



$$\int_{\tau_L}^{\tau_H} e^{-\tau} \tau^{\frac{n}{2}-1} d\tau$$

$$Q^{-\frac{N}{2}}$$

$$Q \;=\; N \left(\mu - \bar{d}\right)^2 + N\, r^2$$

$$\tau_H = \frac{Q}{2 * \sigma_L^2}$$

$$\tau_L = \frac{Q}{2 * \sigma_H^2}$$

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005

# Mean: Unknown noise $\sigma$

provided $\sigma_L \ll r$ and $\sigma_H \gg r$, where $r =$ the RMS residual of the most probable model fit,

**Answer**

$$p(\mu|D,I) \approx \frac{\left[1 + \frac{(\mu - \bar{d})^2}{r^2}\right]^{-\frac{N}{2}}}{\int_{\mu_L}^{\mu_H} d\mu \left[1 + \frac{(\mu - \bar{d})^2}{r^2}\right]^{-\frac{N}{2}}},$$

where the quantity $Nr^2 = \sum (d_i - \bar{d})^2$

Now compare

$$\left[ 1 + \frac{(\mu - \overline{d})^2}{r^2} \right]^{-\frac{N}{2}} \qquad (9.36)$$

with the Student's $t$ distribution

$$f(t|\nu) = \frac{\Gamma[\frac{(\nu+1)}{2}]}{\sqrt{\pi\nu}\Gamma(\frac{\nu}{2})} \left[ 1 + \frac{t^2}{\nu} \right]^{-\frac{(\nu+1)}{2}} .$$

If we set

$$\frac{t^2}{\nu} = \frac{(\mu - \overline{d})^2}{r^2},$$

and the number of degrees of freedom $\nu = N - 1$, then equation (9.36) has the same form as the Student's $t$ distribution .
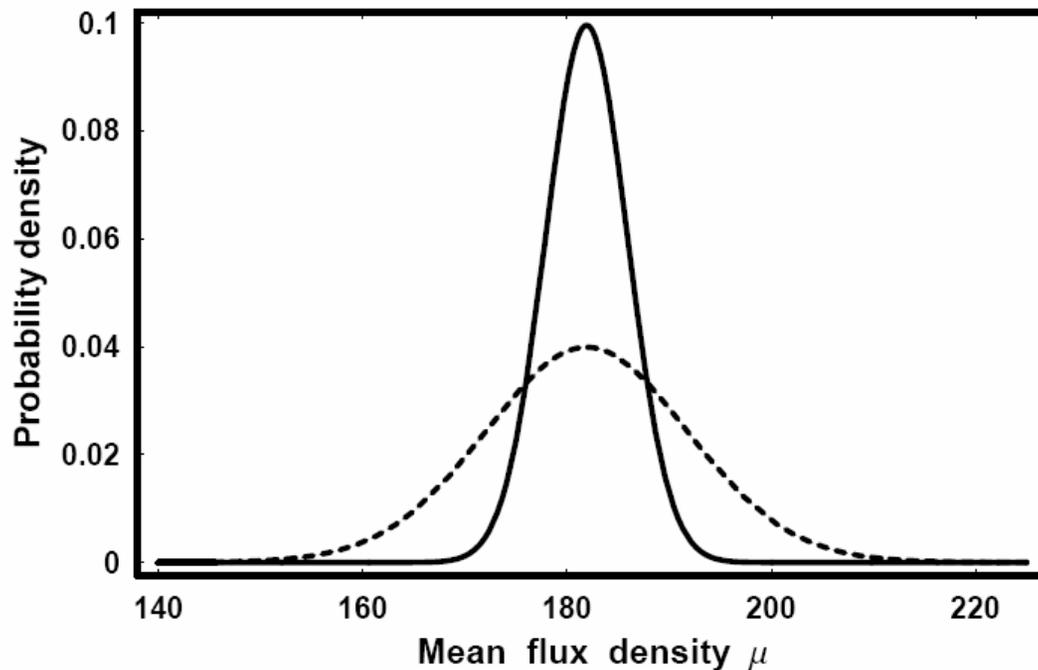
**Fig. Comparison of the computed results for the posterior PDF for the mean radio flux density assuming (a) σ known (solid curve), and (b) marginalizing over an unknown σ (dashed curve).**
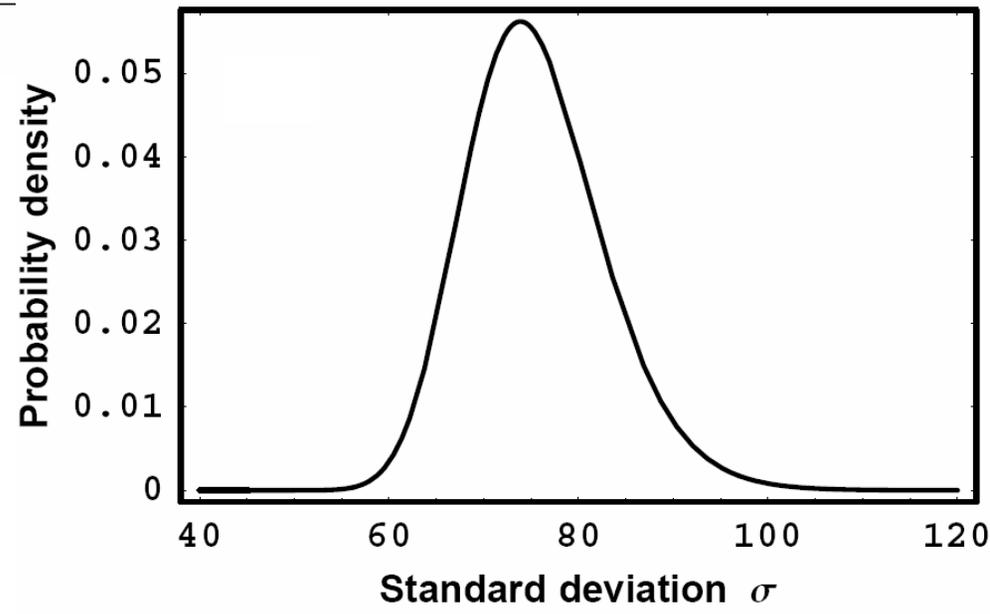
**Marginalizing over σ, in effect estimates σ from the data. Any variability which is not described by the model is assumed to be noise.**

**This leads to a broader posterior, p($\mu$ | D,I), which reflects the larger effective noise.**

## Bayesian estimate of $\sigma$

$$p(\sigma|D,I) = \frac{p(\sigma|I) \int p(\mu|I)p(D|\mu,\sigma,I)d\mu}{p(D|I)}$$

**The distribution is not symmetric like a Gaussian.**

$(\hat{\sigma}, \langle\sigma\rangle, \sqrt{\langle\sigma^2\rangle})$ of $p(\sigma|D,I)$

**are all different.**

$$\hat{\sigma} = r$$

$r$ **= the RMS deviation from** $\overline{d}$.

**and**

$$\langle\sigma^2\rangle = \frac{1}{N-1}\sum_{i=1}^{N}(d_i - \overline{d})^2$$

**Compare the later with the frequentist sample variance,**

$$S^2 = \sum(d_i - \overline{d})^2/(N-1)$$



RMS residual $r = 2$

*Bayesian Logical Data Analysis for the Physical Sciences* © Cambridge University Press 2005