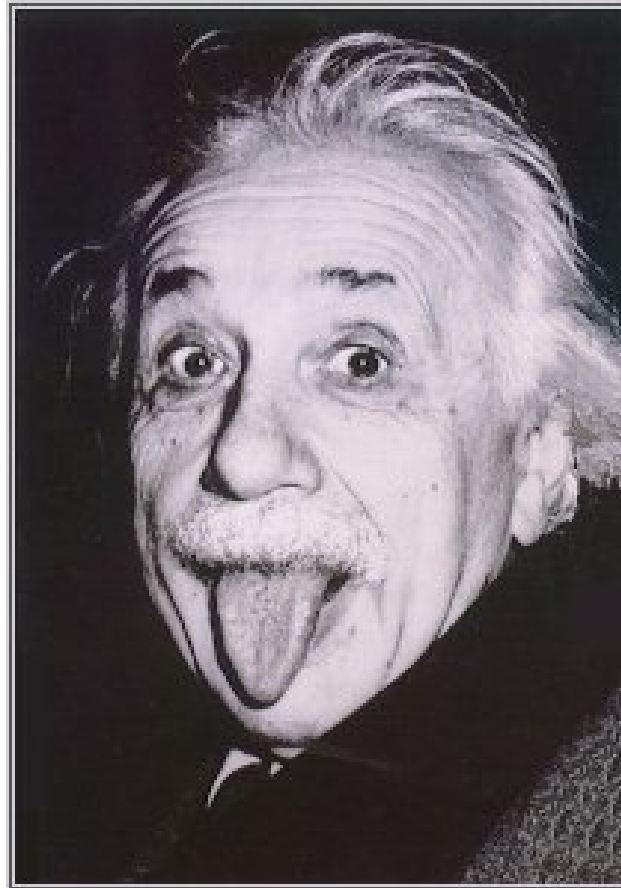


5. Bayesian Model Selection





“Everything should be made as simple as possible, but not simpler”

Bayes' theorem:

$$p(Y | X) = \frac{p(X | Y) \times p(Y)}{p(X)}$$

Laplace rediscovered work of
Rev. Thomas Bayes (1763)



Thomas Bayes
(1702 – 1761 AD)

Posterior

Likelihood

Prior

$$p(\theta \mid \text{data}, M) = \frac{p(\text{data} \mid \theta, M) \times p(\theta \mid M)}{p(\text{data} \mid M)}$$

Evidence



Thomas Bayes
(1702 – 1761 AD)

$$\begin{array}{c}
 \text{Posterior} \\
 \downarrow \\
 p(\theta \mid \text{data}, M) = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}} \\
 \downarrow \qquad \qquad \qquad \downarrow \\
 p(\text{data} \mid \theta, M) \times p(\theta \mid M) \\
 \downarrow \\
 p(\text{data} \mid M)
 \end{array}$$

$$\text{Evidence} = \int p(\text{data} \mid \theta, M) p(\theta \mid M) d\theta$$

Average likelihood, weighted by prior

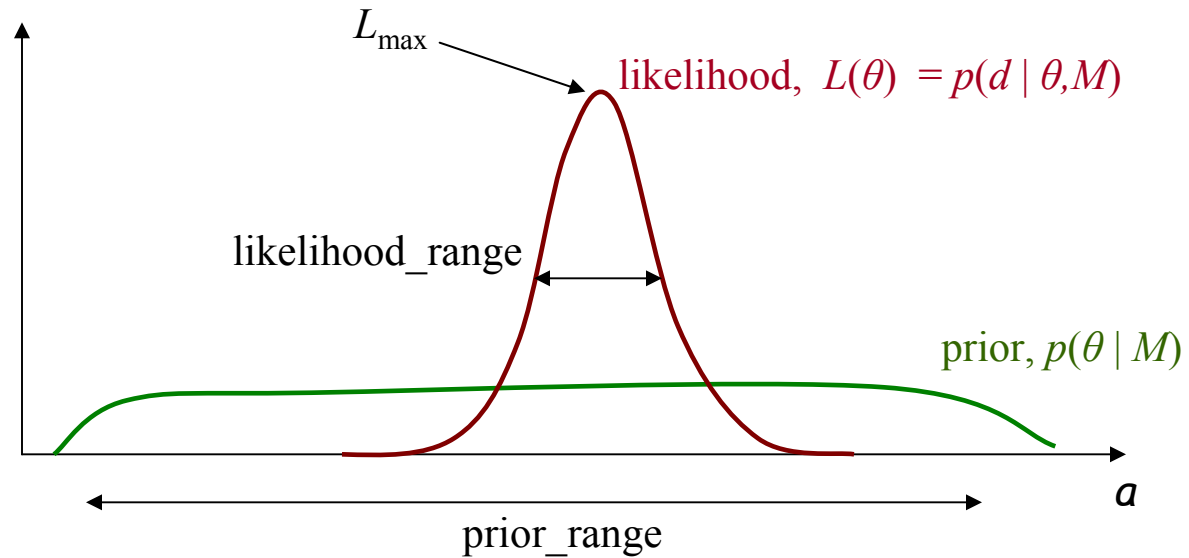
Selecting Between Competing Models

- We can compute the **odds ratio** of two competing models. This can be divided into the **prior odds** and the **Bayes factor**

$$O_{12} = \frac{\text{prob}(M_1 | d)}{\text{prob}(M_2 | d)} = \underbrace{\frac{\text{prob}(M_1)}{\text{prob}(M_2)}}_{\text{prior odds}} \times \underbrace{\frac{\text{prob}(d | M_1)}{\text{prob}(d | M_2)}}_{\text{Bayes factor}}$$

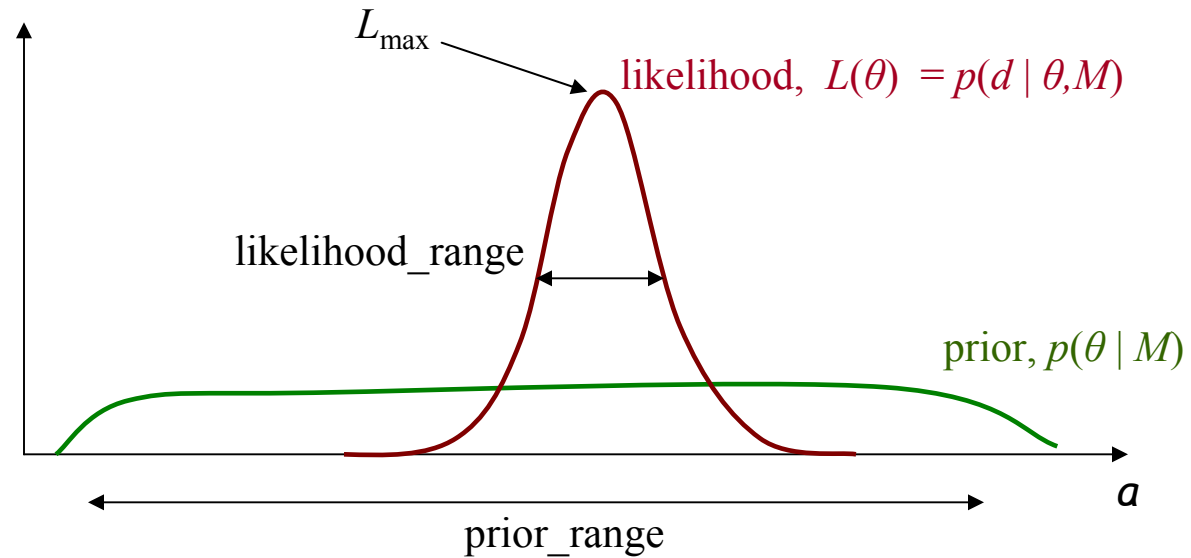
- The Bayes factor is just the ratio of the evidences.

We can split the evidence into two approximate parts:
the maximum of the likelihood and an “Occam factor”:



$$p(d | M) = \int p(\theta | M) p(d | \theta, M) d\theta \approx L_{\max} \underbrace{\frac{\text{likelihood_range}}{\text{prior_range}}}_{\text{the 'Occam factor'}}$$

We can split the evidence into two approximate parts:
the maximum of the likelihood and an “Occam factor”:



$$p(d | M) = \int p(\theta | M) p(d | \theta, M) d\theta \approx L_{\max} \underbrace{\frac{\text{likelihood_range}}{\text{prior_range}}}_{\text{the 'Occam factor'}}$$

The Occam factor penalises models that include wasted parameter space, even if they show a good ML fit.

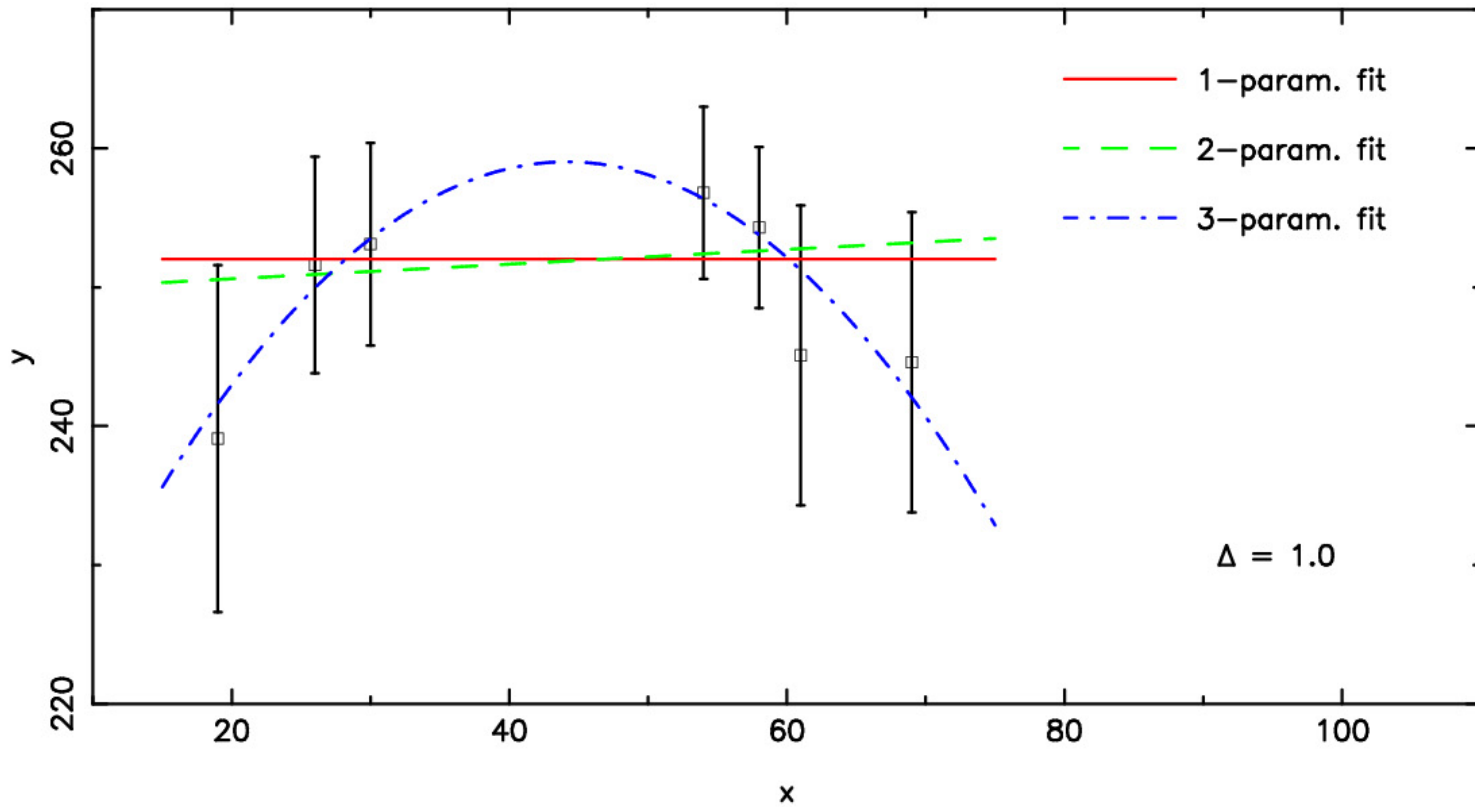


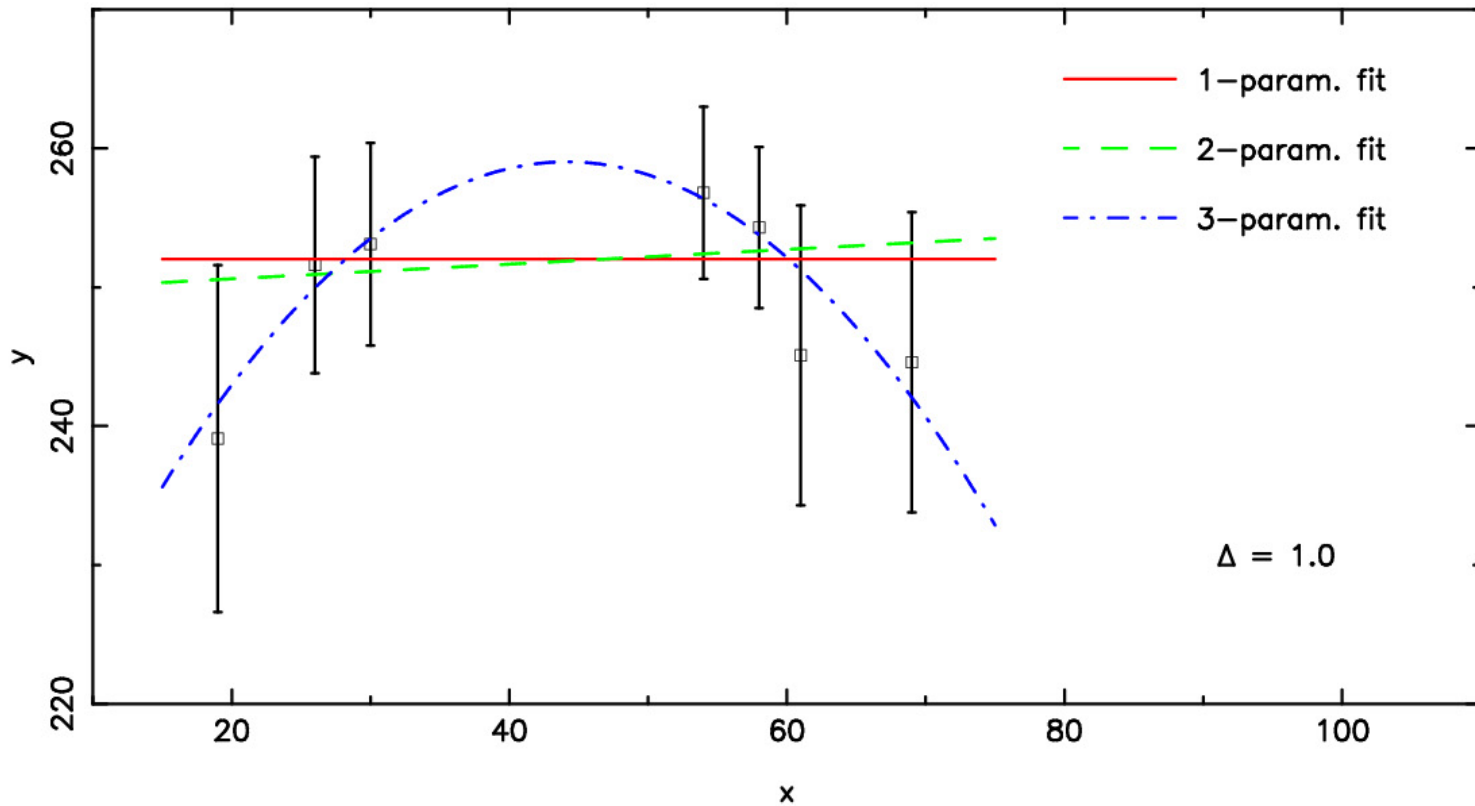
William of Ockham
(1288 – 1348 AD)

Occam's Razor

“It is vain to do with more what can be done with less.”

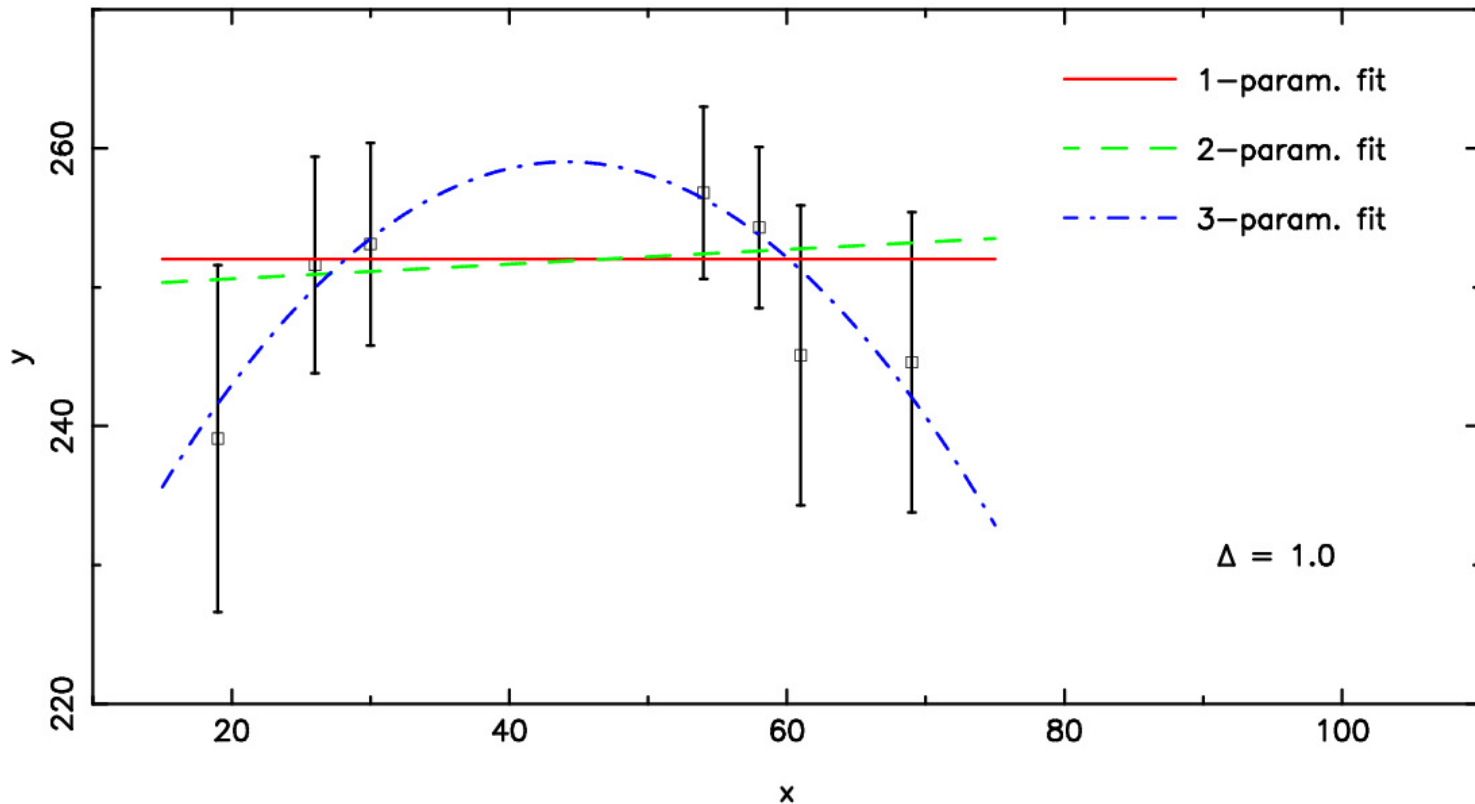
Everything else being equal, we favour models which are ***simple***.





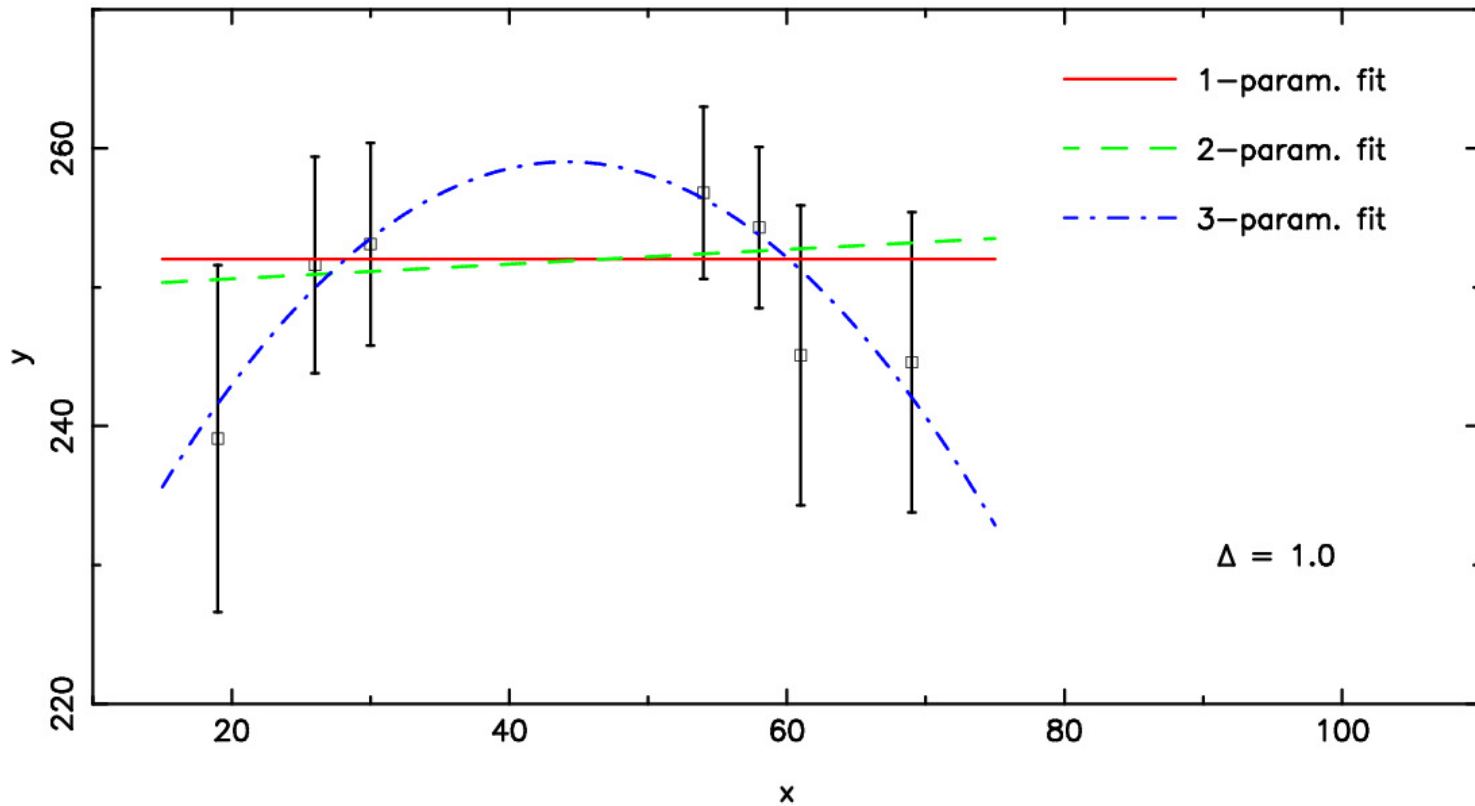
We can compute e.g. $O_{12} = \frac{\text{prob}(m = 1 | d)}{\text{prob}(m = 2 | d)}$

$$O_{13} = \frac{\text{prob}(m = 1 | d)}{\text{prob}(m = 3 | d)}$$



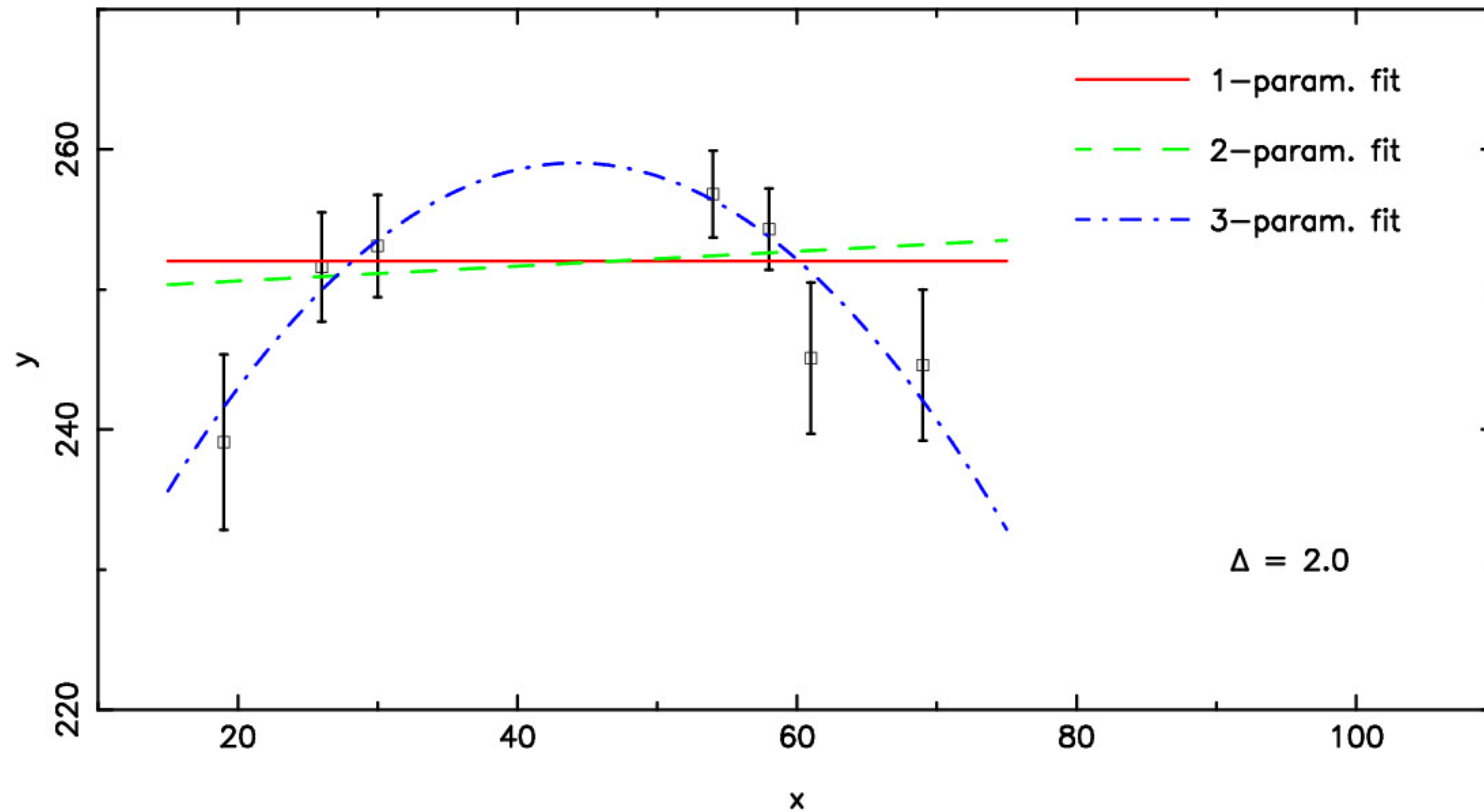
We can compute e.g. $O_{12} = \frac{\text{prob}(m = 1 | d)}{\text{prob}(m = 2 | d)} = 7.2$

$$O_{13} = \frac{\text{prob}(m = 1 | d)}{\text{prob}(m = 3 | d)} = 172.0$$



What if the error bars were over-estimated?

e.g. divide by factor Δ

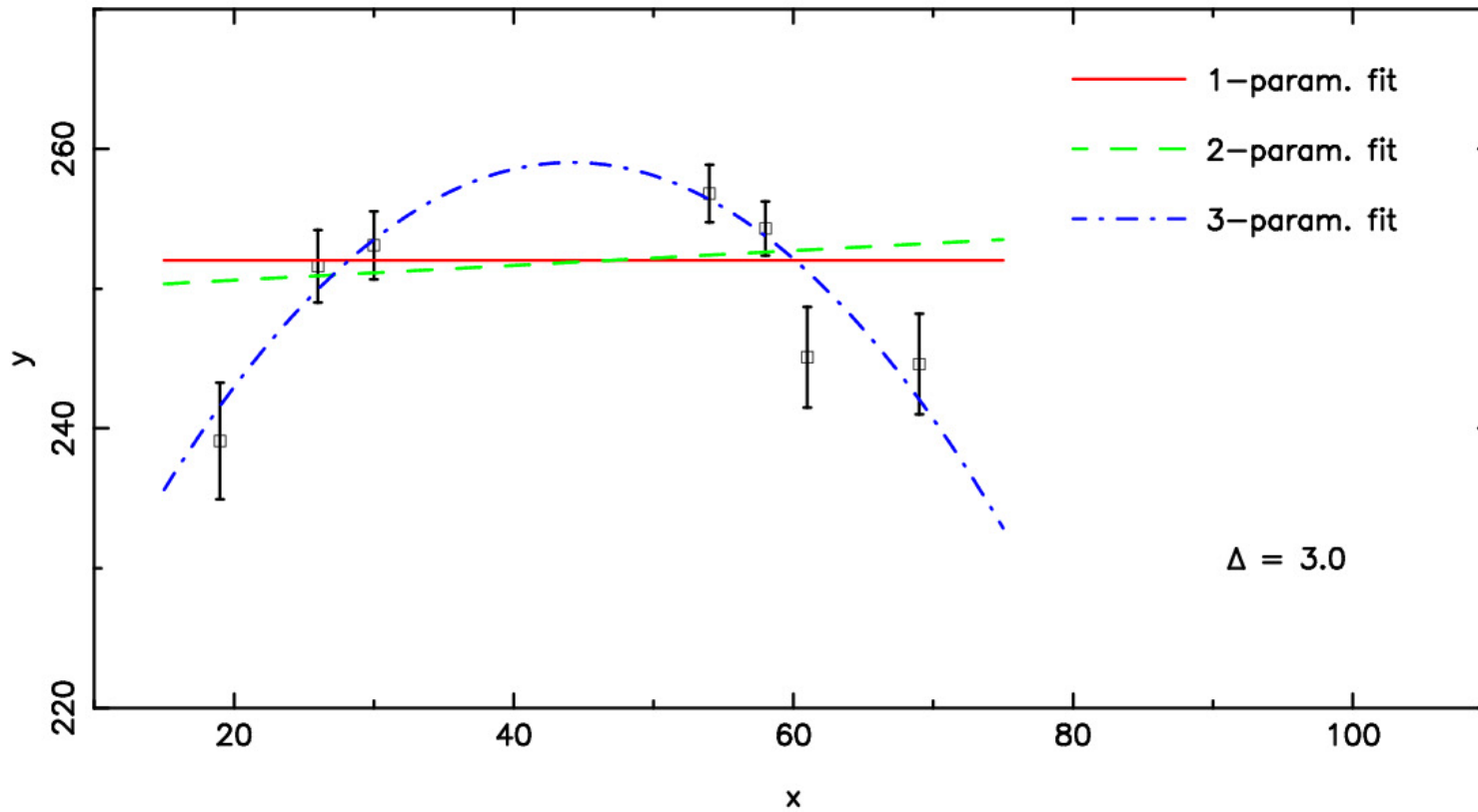


What if the error bars were over-estimated?

e.g. divide by factor $\Delta = 2.0$

$$O_{12} = 6.4$$

$$O_{13} = 5.9$$



What if the error bars were over-estimated?

e.g. divide by factor $\Delta = 3.0$

$$O_{12} = 5.2$$

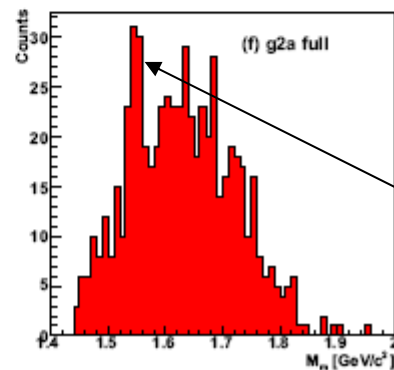
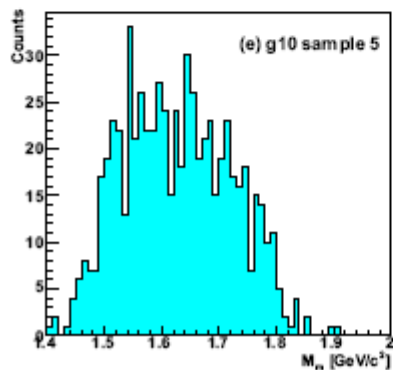
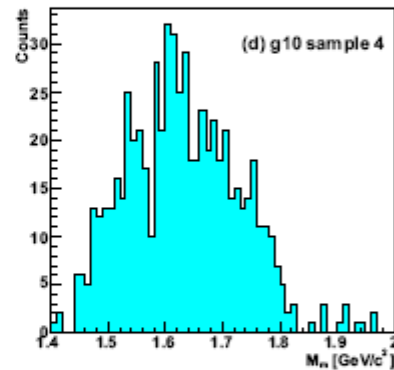
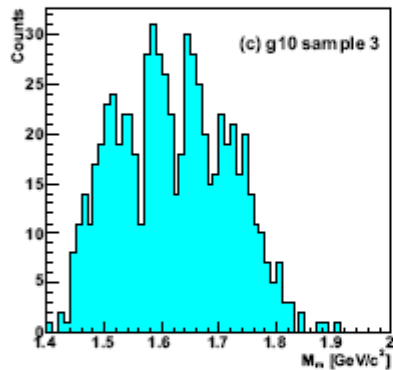
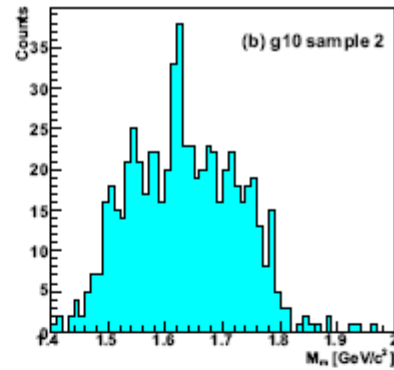
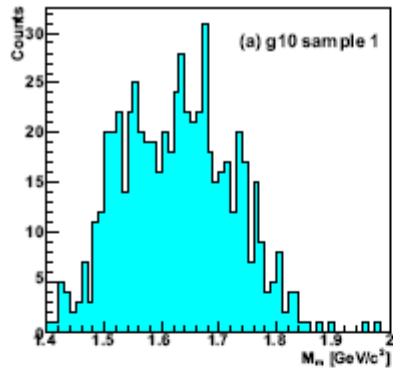
$$O_{13} = 0.02$$

Ireland et al. (2008)

A Bayesian analysis of pentaquark signals from CLAS data

D.G. Ireland,¹ B. McKinnon,¹ D. Protopopescu,¹ P. Ambrozewicz,¹³ M. Anghinolfi,¹⁸ G. Asryan,³⁸ H. Avakian,³³ H. Bagdasaryan,²⁸ N. Baillie,³⁷ J.P. Ball,³ N.A. Baltzell,³² V. Batourine,²² M. Battaglieri,¹⁸ I. Bedlinskiy,²⁰ M. Bellis,⁶ N. Benmouna,¹⁵ B.L. Berman,¹⁵ A.S. Biselli,^{6,12} L. Blaszczyk,¹⁴ S. Bouchigny,¹⁹ S. Boiarinov,³³ R. Bradford,⁶ D. Branford,¹¹ W.J. Briscoe,¹⁵ W.K. Brooks,³³ V.D. Burkert,³³ C. Butuceanu,³⁷ J.R. Calarco,²⁵ S.L. Careccia,²⁸ D.S. Carman,³³ L. Casey,⁷ S. Chen,¹⁴ L. Cheng,⁷ P.L. Cole,¹⁶ P. Collins,³ P. Coltharp,¹⁴ D. Crabb,³⁶ V. Crede,¹⁴ N. Dashyan,³⁸ R. De Masi,^{8,19} R. De Vita,¹⁸ E. De Sanctis,¹⁷ P.V. Degtyarenko,³³ A. Deur,³³ R. Dickson,⁶ C. Djalali,³² G.E. Dodge,²⁸ J. Donnelly,¹ D. Doughty,^{9,33} M. Dugger,³ O.P. Dzyubak,³² K.S. Egiyan,³⁸ L. El Fassi,² L. Elouadrhiri,³³ P. Eugenio,¹⁴ G. Fedotov,²⁴ G. Feldman,¹⁵ A. Fradi,¹⁹ H. Funsten,³⁷ M. Garçon,⁸ G. Gavalian,²⁸ N. Gevorgyan,³⁸ G.P. Gilfoyle,³¹ K.L. Giovanetti,²¹ F.X. Girod,^{8,33} J.T. Goetz,⁴ W. Gohn,¹⁰ A. Gonenc,¹³ R.W. Gothe,³² K.A. Griffioen,³⁷ M. Guidal,¹⁹ N. Guler,²⁸ L. Guo,³³ V. Gyurjyan,³³ K. Hafidi,² H. Hakobyan,³⁸ C. Hanretty,¹⁴ N. Hassall,¹ F.W. Hersman,²⁵ I. Hleiqawi,²⁷ M. Holtrop,²⁵ C.E. Hyde-Wright,²⁸ Y. Ilieva,¹⁵ B.S. Ishkhanov,²⁴ E.L. Isupov,²⁴ D. Jenkins,³⁵ H.S. Jo,¹⁹ J.R. Johnstone,¹ K. Joo,¹⁰ H.G. Juengst,²⁸ N. Kalantarians,²⁸ J.D. Kellie,¹ M. Khandaker,²⁶ W. Kim,²² A. Klein,²⁸ F.J. Klein,⁷ M. Kossov,²⁰ Z. Krahn,⁶ L.H. Kramer,^{13,33} V. Kubarovsky,^{33,29} J. Kuhn,⁶ S.V. Kuleshov,²⁰ V. Kuznetsov,²² J. Lachniet,²⁸ J.M. Laget,³³ J. Langheinrich,³² D. Lawrence,²³ K. Livingston,¹ H.Y. Lu,³² M. MacCormick,¹⁹ N. Markov,¹⁰ P. Mattione,³⁰ B.A. Mecking,³³ M.D. Mestayer,³³ C.A. Meyer,⁶ T. Mibe,²⁷ K. Mikhailov,²⁰ M. Mirazita,¹⁷ R. Miskimen,²³ V. Mokeev,^{24,33} B. Moreno,¹⁹ K. Moriya,⁶ S.A. Morrow,^{8,19} M. Moteabbed,¹³ E. Munevar,¹⁵ G.S. Mutchler,³⁰ P. Nadel-Turonski,¹⁵ R. Nasseripour,³² S. Niccolai,¹⁹ G. Niculescu,²¹ I. Niculescu,²¹ B.B. Niczyporuk,³³ M.R. Niroula,²⁸ R.A. Niyazov,³³ M. Nozar,³³ M. Osipenko,^{18,24} A.I. Ostrovidov,¹⁴ K. Park,²² E. Pasyuk,³ C. Paterson,¹ S. Anefalos Pereira,¹⁷ J. Pierce,³⁶ N. Pivnyuk,²⁰ O. Pogorelko,²⁰ S. Pozdniakov,²⁰ J.W. Price,⁵ S. Procureur,⁸ Y. Prok,³⁶ B.A. Raue,^{13,33} G. Ricco,¹⁸ M. Ripani,¹⁸ B.G. Ritchie,³ F. Ronchetti,¹⁷ G. Rosner,¹ P. Rossi,¹⁷ F. Sabatié,⁸ J. Salamanca,¹⁶ C. Salgado,²⁶ J.P. Santoro,⁷ V. Sapunenko,³³ R.A. Schumacher,⁶ V.S. Serov,²⁰ Y.G. Sharabian,³³ D. Sharov,²⁴ N.V. Shvedunov,²⁴ L.C. Smith,³⁶ D.I. Sober,⁷ D. Sokhan,¹¹ A. Stavinsky,²⁰ S.S. Stepanyan,²² S. Stepanyan,³³ B.E. Stokes,¹⁴ P. Stoler,²⁹ S. Strauch,³² M. Taiuti,¹⁸ D.J. Tedeschi,³² A. Tkabladze,¹⁵ S. Tkachenko,²⁸ C. Tur,³² M. Ungaro,¹⁰ M.F. Vineyard,³⁴ A.V. Vlassov,²⁰ D.P. Watts,¹¹ L.B. Weinstein,²⁸ D.P. Weygand,³³ M. Williams,⁶ E. Wolin,³³ M.H. Wood,³² A. Yegneswaran,³³ L. Zana,²⁵ J. Zhang,²⁸ B. Zhao,¹⁰ and Z.W. Zhao³²

(The CLAS Collaboration)



- Model M_0 : The spectrum can be described by a 3^{rd} order polynomial in the region of interest. This represents the assumption that there is no new particle. A 3^{rd} order polynomial was employed in the original analysis to model the background shape. This model depends on four parameters.
- Model M_P : The spectrum can be described by a “narrow” Gaussian peak sitting atop a 3^{rd} order polynomial background in the region of interest. “Narrow” in this case meaning that the width is significantly less than the region of interest in the mass spectrum. This model depends on seven parameters.

To compare the different models, a ratio of their probabilities in the light of data can be formed:

$$R_E = \frac{P(M_P | D)}{P(M_0 | D)} = \frac{P(D | M_P)}{P(D | M_0)} \times \frac{P(M_P)}{P(M_0)},$$

Significant peak?

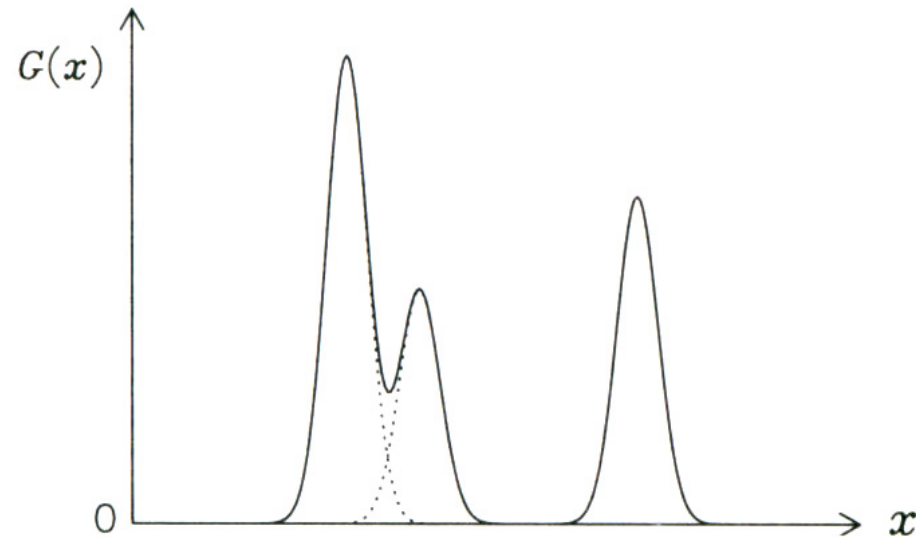
Example from Sivia, Section 4.2: How many spectral lines?

Model: Spectral lines

$$G(x) = \sum_{j=1}^M A_j f(x, x_j),$$

where

$$f(x, x_j) = \exp\left[-\frac{(x - x_j)^2}{2W^2}\right]$$

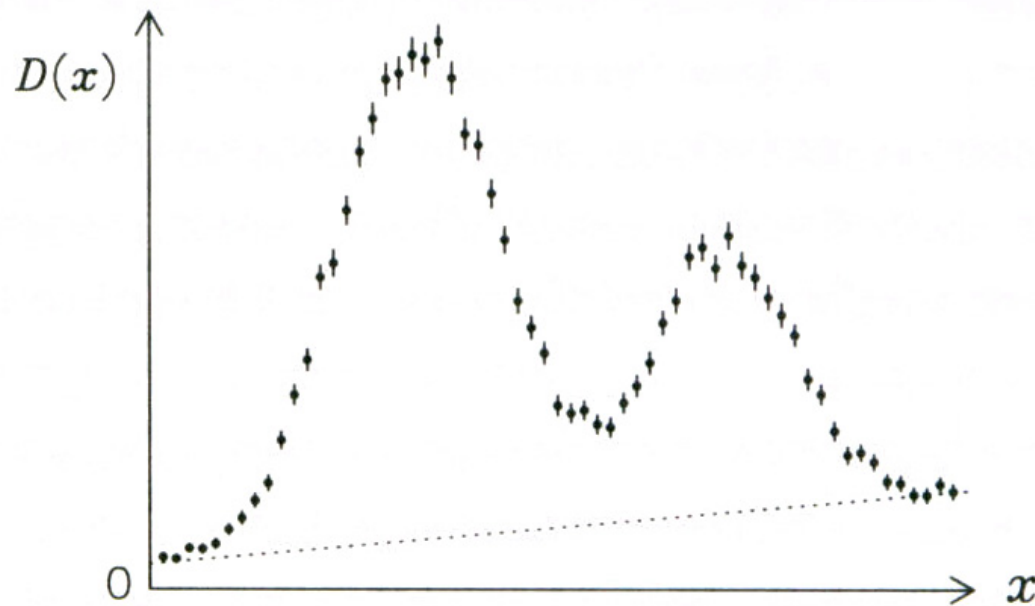


Observed data:

$$D(x_k) = \int G(x) R(x_k - x) dx + B(x_k) + \text{noise}$$

Blurring function
(assumed known)

background



$$\text{prob}(M|\{D_k\}, I) = \frac{\text{prob}(\{D_k\}|M, I) \times \text{prob}(M|I)}{\text{prob}(\{D_k\}|I)}$$

Taking a uniform prior on M implies

$$\text{prob}(M|\{D_k\}, I) \propto \text{prob}(\{D_k\}|M, I)$$

where $\text{prob}(\{D_k\}|M, I) = \iint \cdots \int \text{prob}(\{D_k\}, \{A_j, x_j\}|M, I) d^M A_j d^M x_j$

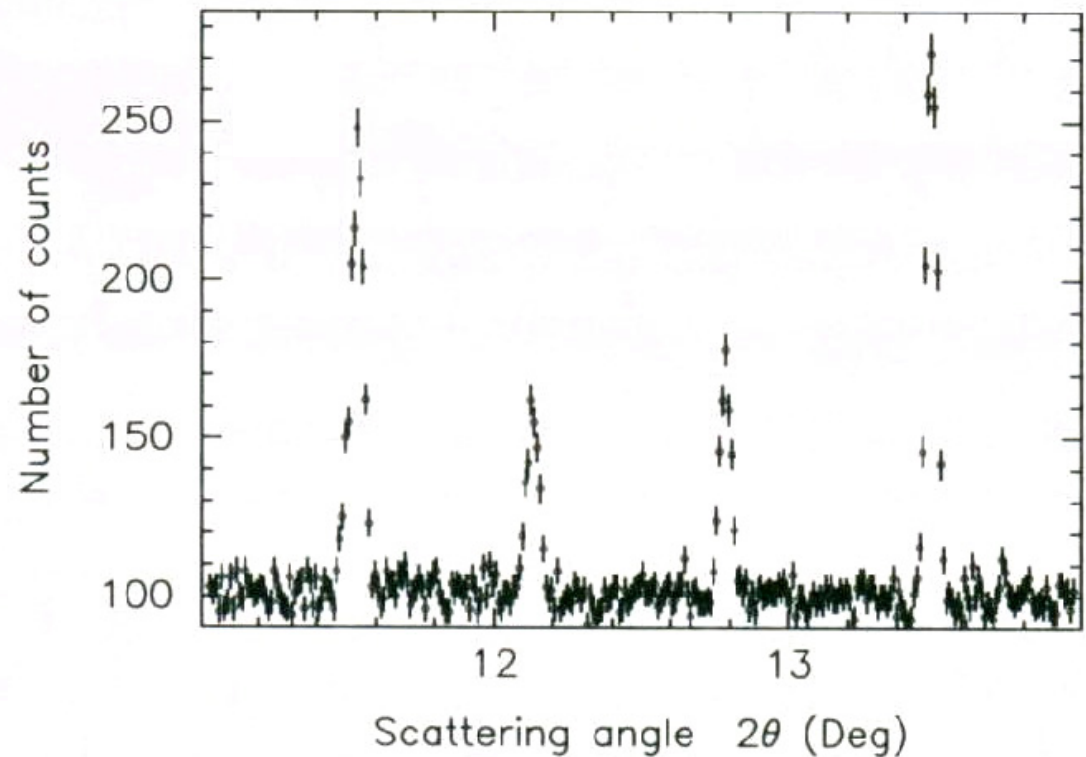
and

$$\text{prob}(\{D_k\}, \{A_j, x_j\}|M, I) = \underbrace{\text{prob}(\{D_k\}|\{A_j, x_j\}, M, I)}_{\text{likelihood}} \underbrace{\text{prob}(\{A_j, x_j\}|M, I)}_{\text{prior}}$$

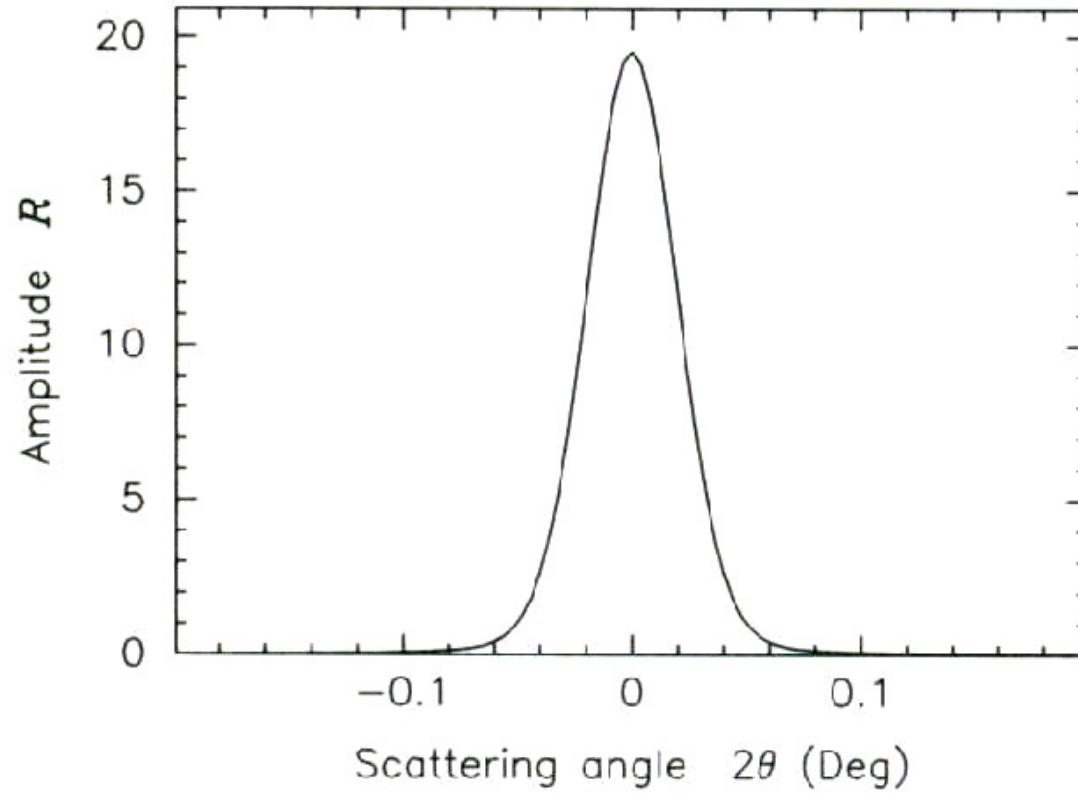
Taking uniform priors on the $\{A_j, x_j\}$ implies

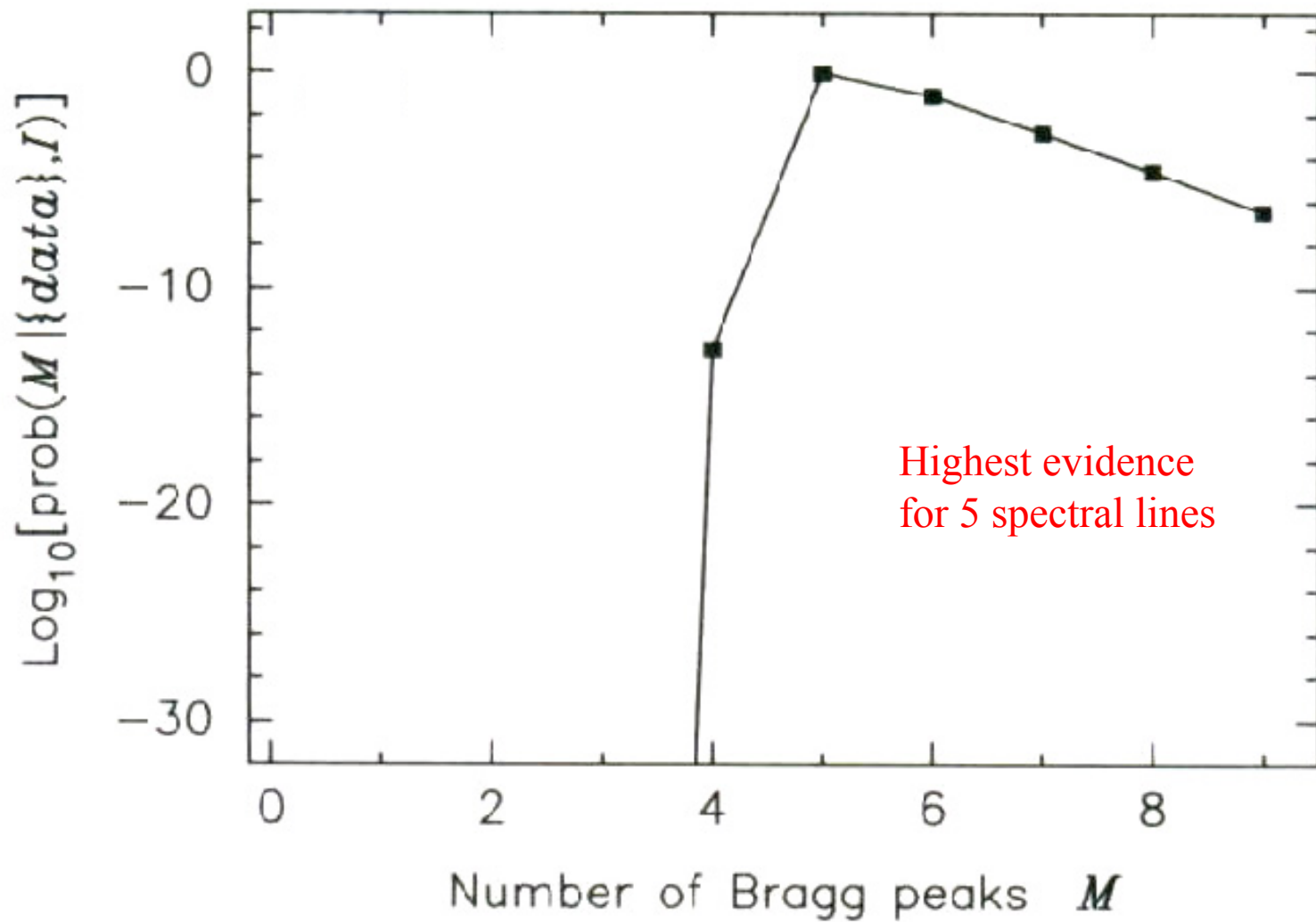
$$\text{prob}(M|\{D_k\}, I) \propto [(x_{\max} - x_{\min}) A_{\max}]^{-M} \iint \cdots \int \exp\left(-\frac{\chi^2}{2}\right) d^M A_j d^M x_j$$

Simulated example



Assume blurring function known....





Question 12: The evidence is smaller for $M > 5$ most likely because

- A** the ML fit is poorer for $M > 5$
- B** the prior on M is smaller for $M > 5$
- C** the improvement in the ML fit for $M > 5$ is more than offset by the reduced Occam factor
- D** none of the above



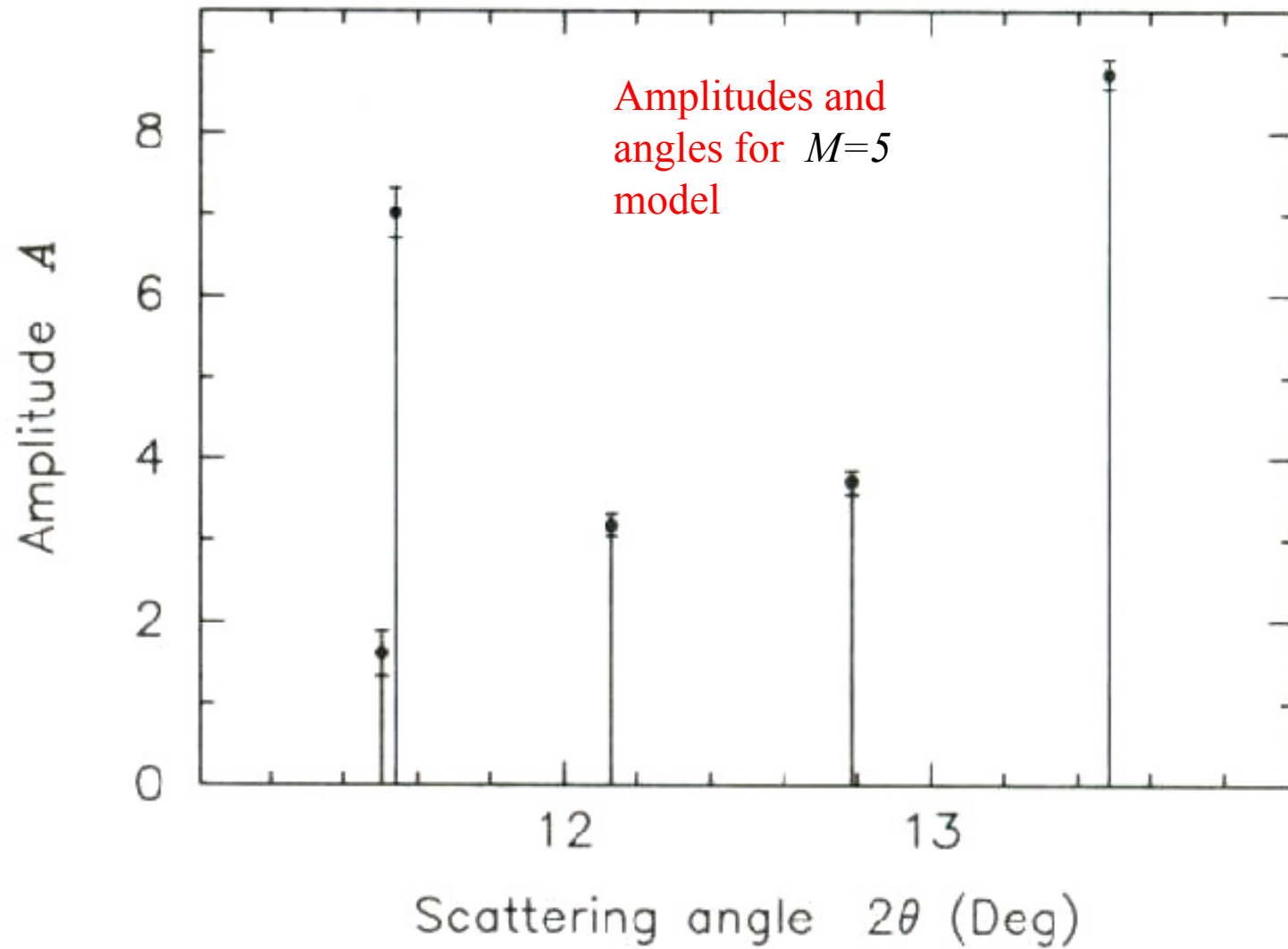
Question 12: The evidence is smaller for $M > 5$ most likely because

A the ML fit is poorer for $M > 5$

B the prior on M is smaller for $M > 5$

C the improvement in the ML fit for $M > 5$ is more than offset by the reduced Occam factor

D none of the above



Taking uniform priors on the $\{A_j, x_j\}$ implies

$$\text{prob}(M|\{D_k\}, I) \propto [(x_{\max} - x_{\min}) A_{\max}]^{-M} \iint \cdots \int \exp\left(-\frac{\chi^2}{2}\right) d^M A_j d^M x_j$$

Evaluating this integral can be a major computational challenge

Approximating the Evidence

$$\text{Evidence} = \int p(\text{data} | \theta, M) p(\theta | M) d\theta$$

Average likelihood, weighted by prior

- Calculating the evidence can be computationally very costly (e.g. CMBR C_ℓ spectrum in cosmology)
- How to proceed?...
 1. Information criteria (Liddle 2004, 2007)
 2. Laplace and Savage-Dickey approximations (Trotta 2005)
 3. Nested sampling (Skilling 2004, 2006; <http://www.inference.phy.cam.ac.uk/bayesys/>) (Mukherjee et al. 2005, 2007; Sivia 2006)

Akaike Information Criterion (Akaike 1974)

$$\text{AIC} = -2 \ln L_{\max} + 2k$$

Number of parameters in model

- Models with too few parameters give poor fit → first term large
- Models with too many parameters penalised by second term
- MC testing (e.g. Kass & Rafferty 1995): can favour models with too many parameters
- ‘dimensionally inconsistent’
- Can give useful upper limit on number of parameters

Bayesian Information Criterion (Schwarz 1978)

$$\text{BIC} = -2 \ln L_{\max} + k \ln N$$

Number of datapoints used in fit

- Approximation to the Bayes factor
- Dimensionally consistent
- If $\text{BIC}(1) - \text{BIC}(2) > 2 \Rightarrow$ positive evidence favouring Model 2
- If $\text{BIC}(1) - \text{BIC}(2) > 6 \Rightarrow$ strong evidence favouring Model 2

(Jeffreys 1961; Mukherjee et al. 1998)

Can we do better than the BIC?

- Laplace approximation to the Bayes factor:
assume posterior well described by a **multivariate Gaussian** around best-fit parameters

Following Trotta (2005)

$$\ln \frac{\bar{\mathcal{P}}(\boldsymbol{\theta} | \mathbf{D}, \mathcal{M})}{\bar{\mathcal{P}}(\boldsymbol{\theta}_* | \mathbf{D}, \mathcal{M})} \approx -\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^T \mathbf{C}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_*)$$

Unnormalised posterior

Best-fit (i.e. ML) parameters

Covariance matrix

Comparing models \mathcal{M}_0 and \mathcal{M}_1 , the Bayes factor B_{01} satisfies

$$\ln B_{01} \approx \mathcal{L}_{01} + \mathcal{C}_{01} + \mathcal{F}_{01},$$

where

$$\mathcal{L}_{01} \equiv \ln \frac{L(\mathbf{D} | \boldsymbol{\theta}_*^{(0)}, \mathcal{M}_0)}{L(\mathbf{D} | \boldsymbol{\theta}_*^{(1)}, \mathcal{M}_1)},$$

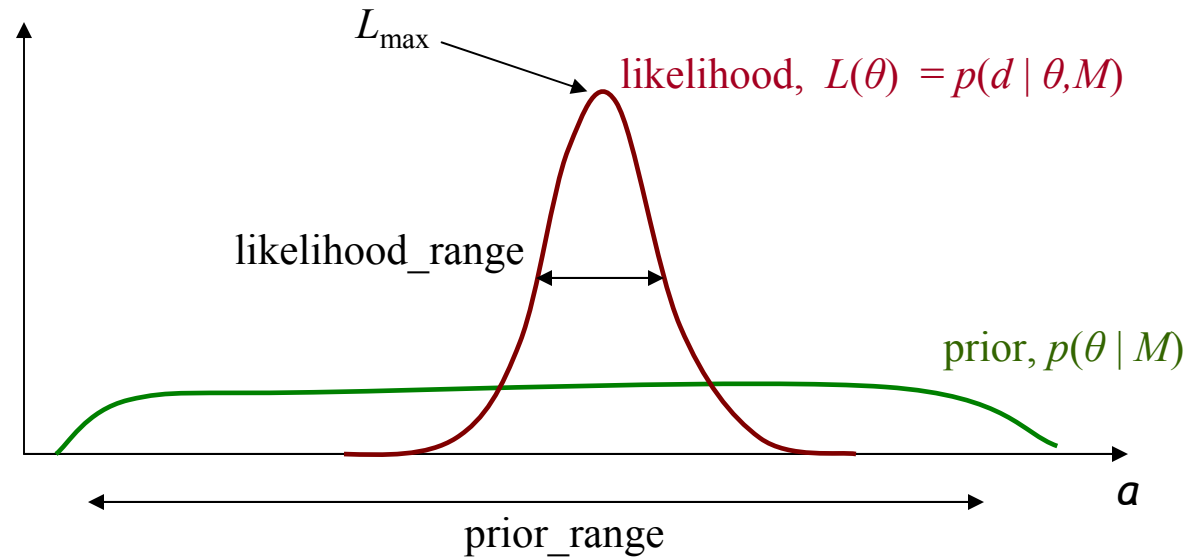
Likelihood ratio

$$\left[\begin{array}{l} \mathcal{C}_{01} \equiv \frac{1}{2} \left(\ln \left[(2\pi)^{d^{(0)} - d^{(1)}} \right] + \ln \frac{\det \mathbf{C}^{(0)}}{\det \mathbf{C}^{(1)}} \right), \\ \mathcal{F}_{01} \equiv \ln \frac{\Delta \boldsymbol{\theta}^{(1)}}{\Delta \boldsymbol{\theta}^{(0)}} \end{array} \right.$$

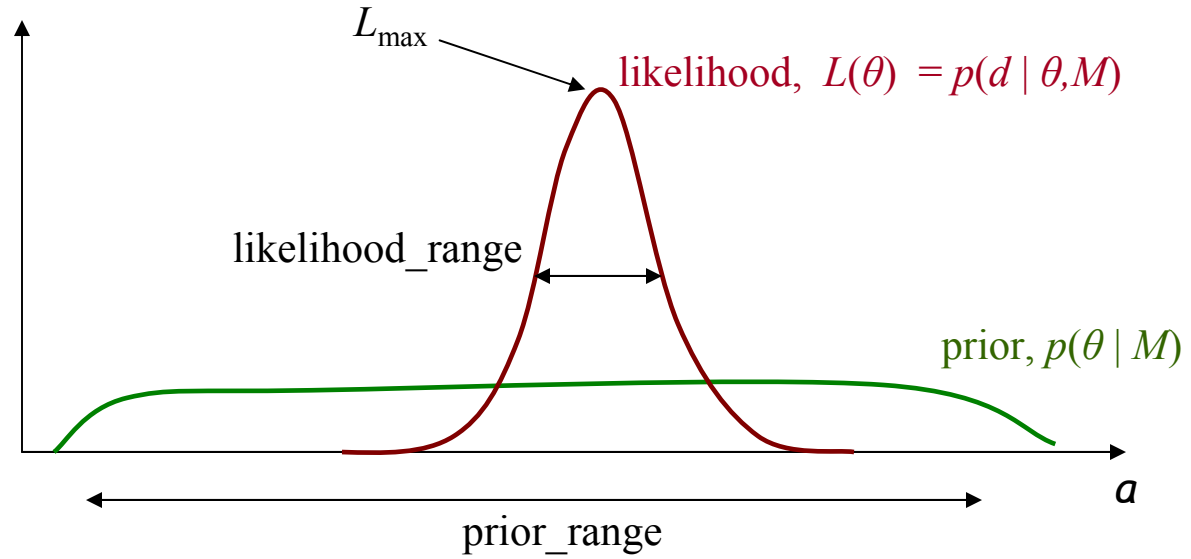
Occam factor

Number of parameters

'Width' of prior



$$p(d | M) = \int p(\theta | M) p(d | \theta, M) d\theta \approx L_{\max} \underbrace{\frac{\text{likelihood_range}}{\text{prior_range}}}_{\text{the 'Occam factor'}}$$



$$p(d | M) = \int p(\theta | M) p(d | \theta, M) d\theta \approx L_{\max} \frac{\text{likelihood_range}}{\text{prior_range}}$$

the 'Occam factor'

\mathcal{L}_{01} \mathcal{F}_{01} \mathcal{C}_{01}

Comparing models \mathcal{M}_0 and \mathcal{M}_1 , the Bayes factor B_{01} satisfies

$$\ln B_{01} \approx \mathcal{L}_{01} + \mathcal{C}_{01} + \mathcal{F}_{01},$$

where

$$\mathcal{L}_{01} \equiv \ln \frac{L(\mathbf{D} | \boldsymbol{\theta}_*^{(0)}, \mathcal{M}_0)}{L(\mathbf{D} | \boldsymbol{\theta}_*^{(1)}, \mathcal{M}_1)},$$

Likelihood ratio

$$\mathcal{C}_{01} \equiv \frac{1}{2} \left(\ln \left[(2\pi)^{d^{(0)} - d^{(1)}} \right] + \ln \frac{\det \mathbf{C}^{(0)}}{\det \mathbf{C}^{(1)}} \right),$$

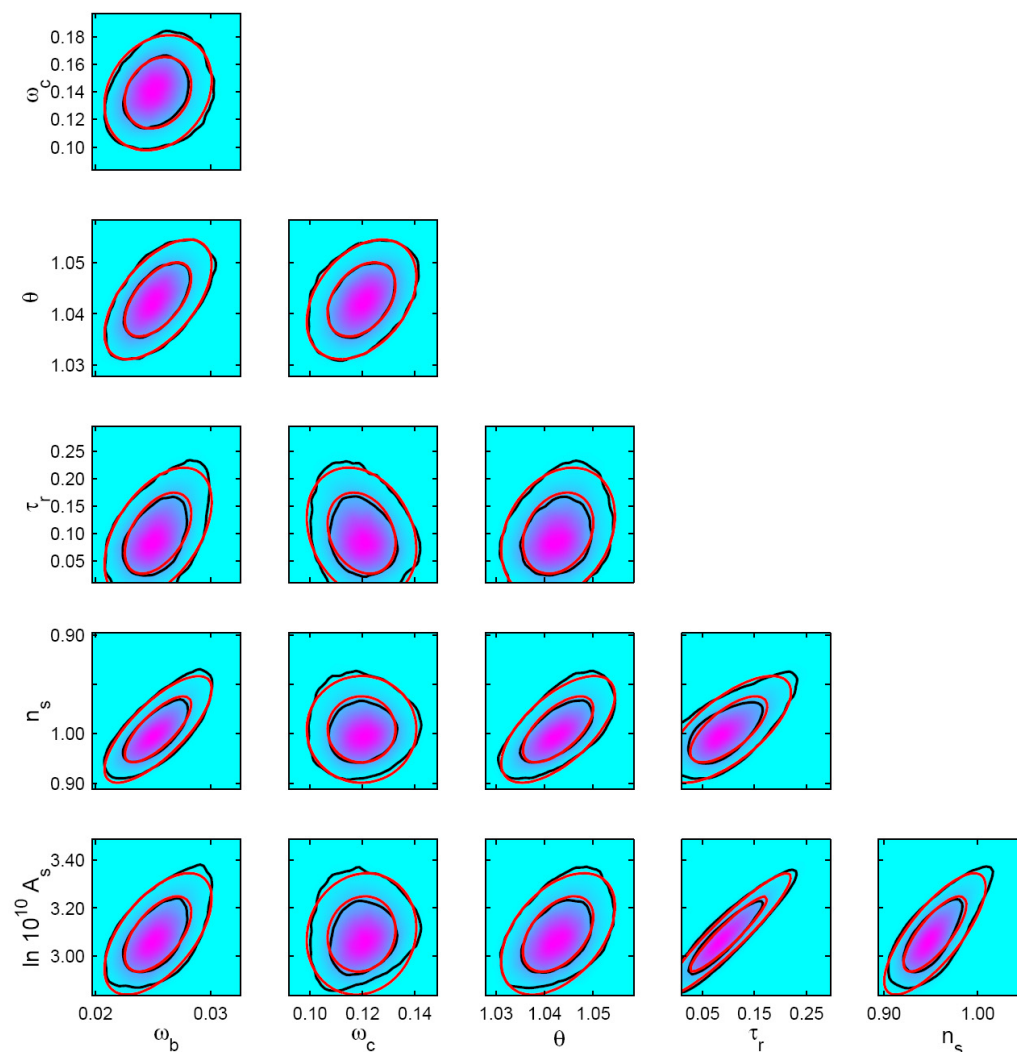
Occam factor

$$\mathcal{F}_{01} \equiv \ln \frac{\Delta \boldsymbol{\theta}^{(1)}}{\Delta \boldsymbol{\theta}^{(0)}}$$

Number of parameters

'Width' of prior

Testing the Laplace approximation



From Trotta (2005)

Good agreement between (MCMC sampled) posteriors and Laplace approximation.

Nested models? Can do better!

Consider \mathcal{M}_1 with two parameters (ω, ψ) ,
and 'submodel' \mathcal{M}_0 with (ω_0, ψ) where $\omega_0 = \text{const.}$

- Assume **separable priors** on parameters.
- Can show (see Dickey 1971; Trotta 2005) that

Savage-Dickey density ratio

$$B_{01} = \frac{\mathcal{P}(\omega_0 | \mathbf{D})}{\pi_1(\omega_0)}$$

Marginalised posterior, evaluated at $\omega = \omega_0$

Prior on ω , evaluated at $\omega = \omega_0$

Nested models? Can do better!

Consider \mathcal{M}_1 with two parameters (ω, ψ) ,
and 'submodel' \mathcal{M}_0 with (ω_0, ψ) where $\omega_0 = \text{const}$.

- Assume **separable priors** on parameters.
- Can show (see Dickey 1971; Trotta 2005) that

Savage-Dickey density ratio

$$B_{01} = \frac{\mathcal{P}(\omega_0 | \mathbf{D})}{\pi_1(\omega_0)}$$

Marginalised posterior, evaluated at $\omega = \omega_0$

Prior on ω , evaluated at $\omega = \omega_0$

- No assumption of Gaussianity required