

Advanced statistical methods for cosmology and astroparticle physics

Roberto Trotta, Imperial College London



A word about statistics:

90% of the game is half mental.

Yogi Berra



Contents



• 1. Foundational aspects: what is probability?

Probability as frequency; Probability as degree of knowledge; Bayes Theorem; Priors; Building the likelihood function; Combination of multiple observations; Nuisance parameters

• 2. Learning from experience: Bayesian parameter inference

Markov Chain Monte Carlo methods; Importance sampling; Nested sampling; Reporting inferences; Credible regions vs confidence regions; The meaning of sigma

• 3. Bayesian model selection and cosmological applications

The different levels of inference; The Bayesian evidence and the Bayes factor; Computing Bayes factors; Information criteria for approximate model selection; The meaning of significance; Comparison with classical hypothesis testing; Model complexity; Bayesian model averaging

• 4. Experiment optimization and prediction

Fisher matrix formalism; Figures of merit; Expected usefulness of an experiment ; Survey optimization; Experimental utility; Bayesian adaptive exploration.

What you will learn

- What does it mean to say that $\Omega_m = 0.28 \pm 0.02$?
- How do you get plots like this and what do they mean?
- How can you quantitatively compare different models for your observations?





Recommended reading

- R. Trotta, "Bayes in the sky: Bayesian inference and model selection in cosmology" *Contemporary Physics*, 49, 2 (2008), 71-104 (arXiv: 0803.4089)
- Bayesian methods in cosmology, Hobson et al (eds), CUP (2010)
- Tom Loredo's Bayesian papers: <u>http://www.astro.cornell.edu/staff/loredo/bayes/</u> <u>tjl.html</u>
- G. D'Agostini, Probability and Measurement Uncertainty in Physics a Bayesian Primer (1995), hep-ph/9512295
- E.T. Jaynes, Probability Theory: The Logic of Science, CUP (2003)
- D. MacKay, *Information theory, Inference & Learning Algorithms*, CUP (2003) (available for free on the web for onscreen viewing)
- P. Gregory, *Bayesian logical data analysis for the physical sciences*, CUP (2003)

"Why should i bother?"

- Increasingly complex models and data: "chi-square by eye" simply not enough
- "If it's real, better data will show it": but all the action is in the "discovery zone" around 3-4 sigma significance. This is a moving target.
- Don't waste time explaining effects which are not there
- Plan for the future: which is the best strategy? (survey design & optimization)
- In some cases, there will be no better data! (cosmic variance)





Upper 95% limit on neutrino mass as a function of observed value for different statistical methods



Brad Efron (PHYSTAT 2003)



The rise of Bayesian methods in astrophysics

Roberto Trotta 11

Imperial College London



The matter with priors

 In parameter inference, prior dependence will in principle vanish for strongly constraining data.

A sensitivity analysis is mandatory for all Bayesian methods!



Usually our parameter space is multi-dimensional: how should we report inferences for one parameter at the time?



Marginal posterior:

 $P(\theta_1|D) = \int L(\theta_1, \theta_2) p(\theta_1, \theta_2) d\theta_2$

FREQUENTIST

Profile likelihood: $L(\theta_1) = max_{\theta_2}L(\theta_1, \theta_2)$

Roberto Trotta

Imperial College

London

12

Confidence intervals: Frequentist approach

- Likelihood-based methods: determine the best fit parameters by finding the minimum of -2Log(Likelihood) = chi-squared
 - Analytical for Gaussian likelihoods





Credible regions: Bayesian approach

- Use the prior to define a metric on parameter space.
- Bayesian methods: the best-fit has no special status. Focus on region of large posterior probability mass instead.
 - Markov Chain Monte Carlo (MCMC)
 - Nested sampling
 - Hamiltonian MC
- Determine posterior credible regions: e.g. symmetric interval around the mean containing 68% of samples

68% CREDIBLE REGION





The good news

- Imperial College London
- Marginalisation and profiling give exactly identical results for the linear Gaussian case.
- This is not suprising, as we already saw that the answer for the Gaussian case is numerically identical for both approaches
- And now the bad news: THIS IS NOT GENERICALLY TRUE!
- A good example is the **Neyman-Scott problem**:
 - We want to measure the signal amplitude μ_i of N sources with an uncalibrated instrument, whose Gaussian noise level σ is constant but unknown.
 - Ideally, measure the amplitude of calibration sources or measure one source many times, and infer the value of σ

- In the Neyman-Scott problem, no calibration source is available and we can only get 2 measurements per source. So for N sources, we have N+1 parameters and 2N data points.
- The profile likelihood estime of σ converges to a biased value $\sigma/sqrt(2)$ for N $\rightarrow \infty$
- The Bayesian answer has larger variance but is unbiased

Neyman-Scott problem



Tom Loredo, talk at Banff 2010 workshop:







Markov Chain Monte Carlo

20

Exploration with "random scans"

- Points accepted/rejected in a in/out fashion (e.g., 2-sigma cuts)
- No statistical measure attached to density of points: no probabilistic interpretation of results possible, although the temptation cannot be resisted...
- Inefficient in high dimensional parameters spaces (D>5)
- **HIDDEN PROBLEM:** Random scan explore only a very limited portion of the parameter space!

One recent example: Berger et al (0812.0980) pMSSM scans (20 dimensions)



C. F. B Hewett Supers Prejudi [arXi v

Imperial College

London

check to pMSSI

21

Random scans explore only a small fraction of the parameter space

- "Random scans" of a highdimensional parameter space only probe a very limited sub-volume: this is the concentration of measure phenomenon.
- **Statistical fact:** the norm of *D* draws from U[0,1] concentrates around (D/3)^{1/2} with constant variance





Geometry in high-D spaces

• **Geometrical fact:** in *D* dimensions, most of the volume is near the boundary. The volume inside the spherical core of *D*-dimensional cube is negligible.

Together, these two facts mean that random scan only explore a very small fraction of the available parameter space in high-dimesional models.



Key advantages of the Bayesian approach Imperial College

- Efficiency: computational effort scales ~ N rather than k^N as in grid-scanning methods. Orders of magnitude improvement over grid-scanning.
- Marginalisation: integration over hidden dimensions comes for free.
- Inclusion of nuisance parameters: simply include them in the scan and marginalise over them.
- Pdf's for derived quantities: probabilities distributions can be derived for any function of the input variables

$P(\theta|d, I) \propto P(d|\theta, I)P(\theta|I)$

- Once the RHS is defined, how do we evaluate the LHS?
- Analytical solutions exist only for the simplest cases (e.g. Gaussian linear model)
- Cheap computing power means that numerical solutions are often just a few clicks away!
- Workhorse of Bayesian inference: Markov Chain Monte Carlo (MCMC) methods. A procedure to generate a list of samples from the posterior.

$P(\theta|d, I) \propto P(d|\theta, I) P(\theta|I)$

- A Markov Chain is a list of samples θ₁, θ₂, θ₃,... whose density reflects the (unnormalized) value of the posterior
- A MC is a sequence of random variables whose (n+1)-th elements only depends on the value of the n-th element
- Crucial property: a Markov Chain converges to a stationary distribution, i.e. one that does not change with time. In our case, the posterior.
- From the chain, expectation values wrt the posterior are obtained very simply:

$$\langle \theta \rangle = \int d\theta P(\theta | d) \theta \approx \frac{1}{N} \sum_{i} \theta_{i}$$
$$\langle f(\theta) \rangle = \int d\theta P(\theta | d) f(\theta) \approx \frac{1}{N} \sum_{i} f(\theta_{i})$$

- Once $P(\theta|d, I)$ found, we can report inference by:
 - Summary statistics (best fit point, average, mode)
 - Credible regions (e.g. shortest interval containing 68% of the posterior probability for θ). Warning: this has **not** the same meaning as a frequentist confidence interval! (Although the 2 might be formally identical)
 - Plots of the marginalised distribution, integrating out nuisance parameters (i.e. parameters we are not interested in). This generalizes the propagation of errors:

$$P(\theta|d,I) = \int d\phi P(\theta,\phi|d,I)$$



$$P(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right)$$

Notation:
$$x \sim N(\mu, \sigma^2)$$

• Frequentist statistics (Fisher, Neymann, Pearson):

E.g., estimation of the mean μ of a Gaussian distribution from a list of observed samples x_1, x_2, x_3 ...

The sample mean is the Maximum Likelihood estimator for μ :

 $\mu_{ML} = X_{av} = (x_1 + x_2 + x_3 + \dots x_N)/N$

• Key point:

in P(X_{av}), X_{av} is a random variable, i.e. one that takes on different values across an ensemble of infinite (imaginary) identical experiments. X_{av} is distributed according to X_{av} ~ N(μ , σ^2 /N) for a fixed true μ

The distribution applies to imaginary replications of data.



• Frequentist statistics (Fisher, Neymann, Pearson): The final result for the confidence interval for the mean

$$P(\mu_{ML} - \sigma/N^{1/2} < \mu < \mu_{ML} + \sigma/N^{1/2}) = 0.683$$

• This means:

If we were to repeat this measurements many times, and obtain a 1-sigma distribution for the mean, the true value μ would lie inside the so-obtained intervals 68.3% of the time

 This is not the same as saying: "The probability of µ to lie within a given interval is 68.3%". This statement only follows from using Bayes theorem.

• Bayesian statistics (Laplace, Gauss, Bayes, Bernouilli, Jaynes):

After applying Bayes therorem $P(\mu | X_{av})$ describes the distribution of our degree of belief about the value of μ given the information at hand, i.e. the observed data.

- Inference is conditional only on the observed values of the data.
- There is no concept of repetition of the experiment.

Gaussian case





Non-Gaussian posteriors

Imperial College London



MCMC estimation

- Marginalisation becomes trivial: create bins along the dimension of interest and simply count samples falling within each bins ignoring all other coordinates
- Examples (from **superbayes.org**) :



Fancier stuff







- Several (sophisticated) algorithms to build a MC are available: e.g. Metropolis-Hastings, Hamiltonian sampling, Gibbs sampling, rejection sampling, mixture sampling, slice sampling and more...
- Arguably the simplest algorithm is the **Metropolis (1954) algorithm:**
 - pick a starting location θ_0 in parameter space, compute $P_0 = p(\theta_0|d)$
 - pick a candidate new location θ_c according to a proposal density $q(\theta_0, \theta_1)$
 - evaluate $P_c = p(\theta_c | d)$ and accept θ_c with probability $\alpha = \min\left(\frac{P_c}{P_0}, 1\right)$
 - if the candidate is accepted, add it to the chain and move there; otherwise stay at θ_0 and count this point once more.

Imperial College

London

- Except for simple problems, achieving good MCMC convergence (i.e., sampling from the target) and mixing (i.e., all chains are seeing the whole of parameter space) can be tricky
- There are several diagnostics criteria around but none is fail-safe. Successful MCMC remains a bit of a black art!
- Things to watch out for:
 - Burn in time
 - Mixing
 - Samples auto-correlation

MCMC diagnostics

(d)url

10*

Steps

104

108 Steps Imperial College London

Burn in Mixing Power spectrum 10⁰ 1000 1000 731 10 10 722 L (N) 0.1 0.2 0.3 Ω_h² 0.03 Ω_sh^e 60 80 100 1000 000 10-108 (d)u1 -102 10 10⁻³ 10^{-2} 10^{-1} 722 10 10 k m_{1/2} (GeV)

00 200 300 A_s[10⁻⁵]

100

3 4 A[10⁻⁵]

5

20 z_{re}

30

2

10

(see astro-ph/0405462 for details)

Roberto Trotta

 10^{0}
Bayesian model comparison

Bayesian inference chain

- Select a model (parameters + priors)
- Compute observable quantities as a function of parameters
- Compare with available data
 - derive parameters constraints: **PARAMETER INFERENCE**
 - compute relative model probability: **MODEL COMPARISON**
- Go back and start again

Imperial College

London

The 3 levels of inference

Imperial College London

LEVEL 1 I have selected a model M and prior P(**θ**|M)

LEVEL 2

Actually, there are several possible models: M₀, M₁,...

LEVEL 3

None of the models is clearly the best



 $P(\theta|d, M) = \frac{P(d|\theta, M)P(\theta|M)}{P(d|M)}$

Parameter inference

(assumes M is the true model)

odds =
$$\frac{P(M_0|d)}{P(M_1|d)}$$

Model comparison

What is the relative plausibility of M₀, M₁,... in light of the data?



 $P(\theta|d) = \sum_{i} P(M_i|d) P(\theta|d, M_i)$

Model averaging

What is the inference on the parameters accounting for model uncertainty?



$$P(\theta|d, M) = \frac{P(d|\theta, M)P(\theta|M)}{P(d|M)}$$

Bayesian evidence or model likelihood

The evidence is the integral of the likelihood over the prior:

$$P(d|M) = \int_{\Omega} d\theta P(d|\theta, M) P(\theta|M)$$

Bayes' Theorem delivers the model's posterior:

$$P(M|d) = \frac{P(d|M)P(M)}{P(d)}$$

When we are comparing two models:

$$\frac{P(M_0|d)}{P(M_1|d)} = \frac{P(d|M_0)}{P(d|M_1)} \frac{P(M_0)}{P(M_1)}$$

The Bayes factor:

$$B_{01} \equiv \frac{P(d|M_0)}{P(d|M_1)}$$

Posterior odds = Bayes factor × prior odds

Scale for the strength of evidence

 A (slightly modified) Jeffreys' scale to assess the strength of evidence (Notice: this is empirically calibrated!)

InB	relative odds	favoured model's probability	Interpretation
< 1.0	< 3:1 < 0.750		not worth mentioning
< 2.5	< 12:1	0.923	weak
< 5.0	< 150:1	0.993	moderate
> 5.0	> 150:1	> 0.993	strong

Roberto Trotta

Imperial College

London



- Bayes factor balances quality of fit vs extra model complexity.
- It rewards highly predictive models, penalizing "wasted" parameter space



The evidence as predictive probability

• The evidence can be understood as a function of d to give the predictive probability under the model M:



Roberto Trotta

Imperial College

London

Simple example: nested models

 This happens often in practice: we have a more complex model, M₁ with prior P(θ|M₁), which reduces to a simpler model (M₀) for a certain value of the parameter,

e.g. $\theta = \theta^* = 0$ (nested models)

 Is the extra complexity of M₁ warranted by the data?



Simple example: nested models



The rough guide to model comparison

Imperial College London



evidence:
$$P(d|M) = \int_{\Omega} d\theta P(d|\theta, M) P(\theta|M)$$

Bayes factor: $B_{01} \equiv \frac{P(d|M_0)}{P(d|M_1)}$

- Usually computational demanding: multi-dimensional integral!
- Several techniques available:
 - Brute force: thermodynamic integration
 - Laplace approximation: approximate the likelihood to second order around maximum gives Gaussian integrals (for normal prior). Can be inaccurate.
 - Savage-Dickey density ratio: good for nested models, gives the Bayes factor
 - Nested sampling: clever & efficient, can be used generally

- This methods works for nested models and gives the Bayes factor analytically.
- Assumptions: nested models (M₁ with parameters θ , Ψ reduces to M₀ for e.g. $\Psi = 0$) and separable priors (i.e. the prior P(θ , Ψ |M₁) is uncorrelated with P(θ |M₀))



49

Nested sampling

- Perhaps **the** method to compute the evidence
- At the same time, it delivers samples from the posterior: it is a highly efficient sampler! (much better than MCMC in tricky situations)
- Invented by John Skilling in 2005: the gist is to convert a *n*-dimensional integral in a 1D integral that can be done easily.



Liddle et al (2006)



Nested sampling

(animation courtesy of David Parkinson)

An algorithm originally aimed primarily at the Bayesian evidence computation (Skilling, 2006):

$$X(\lambda) = \int_{\mathcal{L}(\theta) > \lambda} P(\theta) d\theta$$
$$P(d) = \int d\theta \mathcal{L}(\theta) P(\theta) = \int_0^1 X(\lambda) d\lambda$$





The MultiNest algorithm



• Feroz & Hobson (2007)



The egg-box example



• MultiNest reconstruction of the egg-box likelihood:



Ellipsoidal decomposition



Unimodal distribution Multimodal distribution



Courtesy Mike Hobson

Multinest: Efficiency

Imperial College London



Gaussian mixture model:

True evidence: log(E) = -5.27 **Multinest:** Reconstruction: $log(E) = -5.33 \pm 0.11$

Likelihood evaluations ~ 10^4

Thermodynamic integration:

Reconstruction: $log(E) = -5.24 \pm 0.12$ Likelihood evaluations ~ 10^{6}



D	Nlike	efficiency	likes per dimension
2	7000	70%	83
5	18000	51%	7
10	53000	34%	3
20	255000	15%	1.8
30	753000	8%	1.6

Roberto Trotta

Г

^oeak

A "simple" example: how many sources?

Imperial College London



A "simple" example: how many sources?

Imperial College London



A "simple" example: how many sources?

Imperial College London

Feroz and Hobson (2007)



Bayesian reconstruction

7 out of 8 objects correctly identified. Mistake happens because 2 objects very close.



Cluster detection from Sunyaev-Zeldovich effect in cosmic microwave background maps



Feroz et al 2009

Background + 3 point radio sources Background + 3 point radio sources + cluster



Bayesian model comparison: **R = P(cluster | data)/P(no cluster | data)**

 $R = 0.35 \pm 0.05$ $R \sim 10^{33}$

Cluster parameters also recovered (position, temperature, profile, etc)

The cosmological concordance model



G 11	A 37	1.0	D.C	Die	0.1
Competing model	$\Delta N_{\rm p}$	I IN B	Ref	Data	Outcome
Initial conditions Isocurvature modes					
CDM isocurvature + arbitrary correlations Neutrino entropy + arbitrary correlations Neutrino velocity	+1 +4 +1 +4 +1	$ \begin{array}{c} -7.6 \\ -1.0 \\ [-2.5, -6.5]^p \\ -1.0 \\ [-2.5, -6.5]^p \end{array} $	[58] [46] [46] [60]	WMAP3+, LSS WMAP1+, LSS, SN Ia WMAP3+, LSS WMAP1+, LSS, SN Ia WMAP3+, LSS	Strong evidence for adiabaticity Undecided Moderate to strong evidence for adiabaticity Undecided Moderate to strong evidence for adiabaticity
+ arbitrary correlations	+4	-1.0	[46]	WMAP1+, LSS, SN Ia	Undecided
Primordial power spectr	um				
No tilt $(n_s = 1)$	-1	$^{+0.4}_{[-1.1, -0.6]^p}$ -0.7 -0.9 $[-0.7, -1.7]^{p,d}$ -2.0 -2.6 -2.9 $< -3.9^c$	[47] [51] [58] [70] [186] [185] [70] [58] [65]	WMAP1+, LSS WMAP1+, LSS WMAP1+, LSS WMAP3+ WMAP3+, LSS WMAP3+, LSS WMAP3+, LSS WMAP3+, LSS	Undecided Undecided Undecided Undecided $n_s = 1$ weakly disfavoured $n_s = 1$ weakly disfavoured $n_s = 1$ moderately disfavoured $n_s = 1$ moderately disfavoured $m_s = 1$ moderately disfavoured Moderate evidence at best against $n_s \neq 1$
Running	+1	$[-0.6, 1.0]^{p,d}$	[186]	WMAP3+, LSS	No evidence for running
Running of running	1.2	$< 0.2^{\circ}$ $< 0.4^{\circ}$	[166]	WMAP3+, LSS WMAP3+, LSS	Running not required
Large scales cut-off	+2	[1.3, 2.2] ^{p,d}	[186]	WMAP3+, LSS	Weak support for a cut-off
Matter-energy content Non-flat Universe	+1	-3.8	[70]	WMAP3+, HST	Flat Universe moderately favoured
Coupled neutrinos	+1	-3.4 -0.7	[58] [193]	WMAP3+, LSS, HST WMAP3+, LSS	Flat Universe moderately favoured No evidence for non–SM neutrinos
Dark energy sector					
$w(z) = w_{\text{eff}} \neq -1$	+1	$[-1.3, -2.7]^p$ -3.0 -1.1 $[-0.2, -1]^p$ $[-1.6, -2.3]^d$	[187] [50] [51] [188]	SN Ia SN Ia WMAP1+, LSS, SN Ia SN Ia, BAO, WMAP3 SN Ia, GBB	Weak to moderate support for Λ Moderate support for Λ Weak support for Λ Undecided Weak support for Λ
$w(z) = w_0 + w_1 z$	+2	$[-1.5, -3.4]^p$ -6.0	[187] [50]	SN Ia SN Ia	Weak to moderate support for Λ Strong support for Λ
$w(z) = w_0 + w_a \left(1 - a\right)$	+2	$^{-1.8}_{-1.1}$ $[-1.2, -2.6]^d$	[188] [188] [189]	SN Ia, BAO, WMAP3 SN Ia, BAO, WMAP3 SN Ia, GRB	Weak support for Λ Weak support for Λ Weak to moderate support for Λ
Reionization history No reionization ($\tau = 0$) No reionization and no tilt	$^{-1}_{-2}$	$^{-2.6}_{-10.3}$	[70] [70]	WMAP3+, HST WMAP3+, HST	$\tau \neq 0$ moderately favoured Strongly disfavoured

from Trotta (2008)

InB < 0: favours ACDM

- Warning: frequentist hypothesis testing (e.g., likelihood ratio test) cannot be interpreted as a statement about the probability of the hypothesis!
- Example: to test the null hypothesis H₀: θ = 0, draw *n* normally distributed points (with known variance σ²). The χ² is distributed as a chi-square distribution with (*n*-1) degrees of freedom (dof). Pick a significance level α (or p-value, e.g. α = 0.05). If P(χ² > χ²_{obs}) < α reject the null hypothesis.
- This is a statement about the likelihood of observing data as extreme or more extreme than have been measured assuming the null hypothesis is correct.
- It is not a statement about the probability of the null hypothesis itself and cannot be interpreted as such! (or you'll make gross mistakes)
- The use of p-values implies that a hypothesis that may be true can be rejected because it has not predicted observable results that have not actually occurred. (Jeffreys, 1961)

The significance of significance

- Imperial College London
- Important: A 2-sigma result does not wrongly reject the null hypothesis 5% of the time: at least 29% of 2-sigma results are wrong!
 - Take an equal mixture of H₀, H₁
 - Simulate data, perform hypothesis testing for H₀
 - Select results rejecting H₀ at 1-α CL
 - What fraction of those results did actually come from H₀ ("true nulls", should not have been rejected)?

p-value	sigma	fraction of true nulls	lower bound
0.05	1.96	0.51	0.29
0.01	2.58	0.20	0.11
0.001	3.29	0.024	0.018

Recommended: Sellke, Bayarri & Berger, The American Statistician, 55, 1 (2001)

• What if we do not know how to set the prior? For nested models, we can still choose a prior that will maximise the support for the more complex model:



• The absolute upper bound: put all prior mass for the alternative onto the observed maximum likelihood value. Then

$$B < \exp(-\chi^2/2)$$

• More reasonable class of priors: symmetric and unimodal around Ψ =0, then (α = significance level)

$$B < \frac{-1}{\exp(1)\alpha \ln \alpha}$$

If the upper bound is small, no other choice of prior will make the extra parameter significant.

Sellke, Bayarri & Berger, The American Statistician, 55, 1 (2001)

Imperial College London

α	sigma	Absolute bound on InB (B)	"Reasonable" bound on InB (B)
0.05	2.0	2.0 (7:1) <mark>weak</mark>	0.9 (3:1) <mark>undecided</mark>
0.003	3.0	4.5 (90:1) moderate	3.0 (21:1) <mark>moderate</mark>
0.0003	3.6	6.48 (650:1) <mark>strong</mark>	5.0 (150:1) <mark>strong</mark>

A conversion table



p-value	\bar{B}	$\ln \bar{B}$	sigma	category
0.05	2.5	0.9	2.0	
0.04	2.9	1.0	2.1	'weak' at best
0.01	8.0	2.1	2.6	
0.006	12	2.5	2.7	'moderate' at best
0.003	21	3.0	3.0	
0.001	53	4.0	3.3	
0.0003	150	5.0	3.6	'strong' at best
6×10^{-7}	43000	11	5.0	

Rule of thumb:

a n-sigma result should be interpreted as a n-1 sigma result

Application: the spectral tilt

Imperial College London

- Is the spectrum of primordial fluctuations scale-invariant?
- Model comparison:
 n = 1 vs n ≠ 1 (with inflation-motivated prior)
- Results:
 - n ≠ 1 favoured with odds of 17:1 (Trotta 2007)
 - n ≠ 1 favoured with odds of 15:1 (Kunz, Trotta & Parkinson 2007)
 - **n** ≠ 1 favoured with odds of 7:1 (Parkinson 2007 et al 2006)
- Upper bound: odds of 49:1 at best for n ≠ 1 (Gordon and Trotta 2007)

Application: dipole modulation



- Eriksen et al (2004) found hints for a dipolar modulation in WMAP1 ILC map
- Adding a phenomenological dipole pattern improves the chisquare by 9 units (for 3 extra parameters)
- Is this significant evidence?
- Not really: upper bound on B is odds of 9:1 The absolute upper bound is about the same (Gordon and Trotta 2007)



nuuerto Trotta

Model complexity

- "Number of free parameters" is a relative concept. The relevant scale is set by the prior range
- How many parameters can the data support, regardless of whether their detection is significant?
- **Bayesian complexity** or effective number of parameters:

$$C_b = \overline{\chi^2(\theta)} - \chi^2(\widehat{\theta})$$
$$= \sum_i \frac{1}{1 + (\sigma_i / \Sigma_i)^2}$$

Kunz, RT & Parkinson, astro-ph/0602378, Phys. Rev. D 74, 023503 (2006) Following Spiegelhalter et al (2002)

Polynomial fitting

• Data generated from a model with n = 6:



London



How many parameters does the CMB need?



Roberto Trotta

Imperial College

London



$P(\boldsymbol{\theta}|d) = \sum_{i} P(\boldsymbol{\theta}|d, M_{i})P(M_{i}|d)$




- Bayesian model comparison extends parameter inference to the space of models
- The Bayesian evidence (model likelihood) represents the change in the degree of belief in the model after we have seen the data
- Models are rewarded for their predictivity (automatic Occam's razor)
- Prior specification is for model comparison a key ingredient of the model building step. If the prior cannot be meaningfully set, then the physics in the model is probably not good enough.
- Bayesian model complexity can help (together with the Bayesian evidence) in assessing model performance.

Prediction and optimization

The Bayesian perspective

- In the Bayesian framework, we can use present-day knowledge to produce probabilistic forecasts for the outcome of a future measurement
- This is **not** limited to assuming a model/parameter value to be true and to determine future errors
- Many questions of interest today are of model comparison: e.g.
 - is dark energy Lambda or modified gravity?
 - is dark energy evolving with time?
 - Is the Universe flat or not?
 - Is the spectrum of perturbations scale invariant or not?

Predictions for future observations



The predictive distribution

- Use present knowledge (and uncertainty!) to predict what a future measurement will find (with corresponding probability)
- True values: $(\theta_0, \theta_1) = (0, 1)$
- Present-day data: d
- Future data: D

 $P(D|d) = \int d\theta P(D|\theta) P(\theta|d)$

Predictive probability = future likelihood weighted by present posterior



 $y_1 \text{ at } x = -0.5$

Predictive Log Prob $(1-2-3\sigma)$ for measurements at (x1,x2) = (-0.5, -1.0)

Predictive distribution



Extending the power of forecasts

- Thanks to predictive probabilities we can increase the scope and power of forecasts:
- Level 0: assume a model M and a fiducial value for the parameters, θ* produce a forecast for the errors that a future experiment will find *if M and* θ* are the correct choices
- Level 1: average over current parameter uncertainty within M
- Level 2: average over current model uncertainty: replace M by M₁, M₂,...

Roberto Trotta

Imperial College

London

Predictive posterior odds distribution

Bayes factor forecast for Planck



(2006), Parkinson et al (2006)

Experiment design

- The optimization problem is fully specified once we define a utility function U depending on the outcome e of a future observation (e.g., scientific return). We write for the utility U(e, o, θ), where o is the current experiment and θ are the true values of the parameters of interest
- We can then evaluate the **expected utility:**

$$\mathcal{E}[U|e,o] = \int d\theta U(\theta, e, o) P(\theta|o)$$

Example: an astronomer measures y = θ x (with Gaussian noise) at a few points 0 < x < 1. She then has a choice between building 2 equally expensive instruments to perform a new measurement:
1. Instrument (e) is as accurate as today's experiments but extends to much larger values of x (to a maximum x_{max})
2. Instrument (a) is much more accurate but it is built in such a way as has to have a "sweet spot" at a certain value of y, call it y*, and much less accurate elsewhere
Which instrument should she go for?

The answer depends on how good her current knowledge is - i.e. is the current uncertainty on θ* small enough to allow her to target accurately enough x=x* so that she can get to the "sweet spot" y*= θ*x*?
 (try it out for yourself! *Hint: use for the utility the inverse variance of the future posterior*

on θ and assume for the noise levels of experiment (a) the toy model:

$$\tau_a^2 = \tau_*^2 \exp\left(\frac{(y-y_\star)^2}{2\Delta^2}\right)$$

where y^* is the location of the sweet spot and Δ is the width of the sweet spot)

Small uncertainty

Large uncertainty



Making predictions: Dark Energy

Imperial College London

A model comparison question: is dark energy Lambda, i.e. $(w_0, w_a) = (-1, 0)$? How well will the future probe SNAP be able to answer this?

 $\begin{array}{c} 1 \\ 0.5 \\ 0.5 \\ -0.5 \\ -1 \\ -2 \\ -2 \\ -1.5 \\ w_0 \\ \end{array}$

Simulates from LCDM Assumes LCDM is true Ellipse not invariant when changing model assumptios

Fisher Matrix

Bayesian evidence



Simulate from all DE models Assess "model confusion" Allows to discriminate against LCDM



- Predictive distributions incorporate present uncertainty in forecasts for the future scientific return of an experiment
- Experiment optimization requires the specification of an utility function. The "best" experiment is the one that maximises the expected utility.