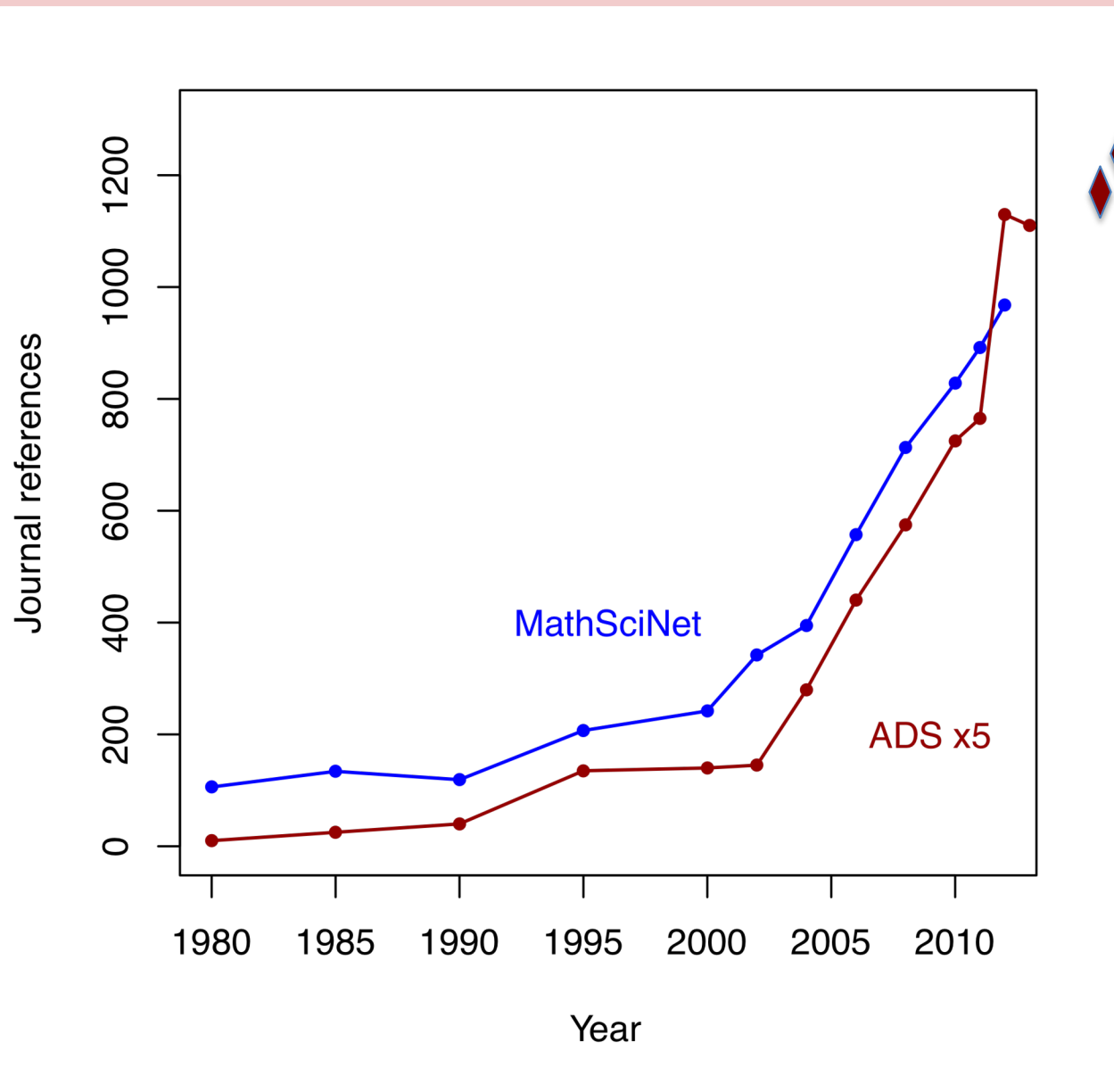


Bayesian inference

Eric Feigelson

2nd East Asian Workshops on Astrostatistics
Nanjing & Guiyang
July 2018

Rapid rise of Bayesian analyses in mathematics and astronomy



Outline

- Derivation of Bayes' Theorem
- A simple astronomical example
- Priors
- Posteriors
- Parameter estimation & credible intervals
- Model selection
- Marginalization
- Hierarchical models
- Philosophical considerations: frequentist or Bayesian?
- Advantages & disadvantages of Bayes
- Final advice

Axioms of probability

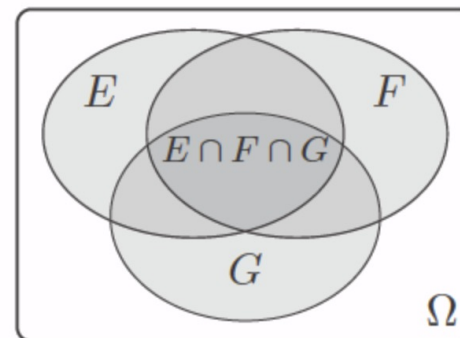
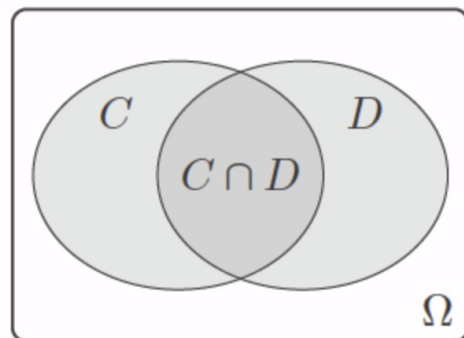
Axiom 1: For an event A , $0 \leq P(A) \leq 1$.

Axiom 2: For a sample space Ω , $P(\Omega) = 1$

Axiom 3: For mutually exclusive events

$$P(A_1 \cup A_2 \cup A_3 \cdots) = P(A_1) + P(A_2) + P(A_3) + \cdots$$

\cup = “or”, union \cap = “and”, intersection $|$ = “given”, conditioned on



Union and intersection of events.

Conditional probability

The probability of event A given event B is equal to the intersection of A and B normalized by the probability of B

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A \cap B) = P(A | B)P(B)$$

Example: The conditional probability that a star has solar mass and a Jovian planet with ellipticity 0.5-0.6 is equal to the product of the probability that a star is a G star (say $P \sim 1\%$) and the probability that any star has a $\epsilon=0.5-0.6$ Jovian planet (say $P \sim 20\%$). The conditional probability is thus $P \sim 0.2\%$.

The 'multiplication rule' easily extends to n events:

$$P(A_1 \cap A_2 \cap \dots A_n) = P(A_1) P(A_2 | A_1) \dots P(A_{n-1} | A_1, \dots A_{n-2}) \\ \times P(A_n | A_1, \dots A_{n-1}).$$

Another astronomical example:

Except for the rare circumstance when an entirely new phenomenon is discovered, astronomers are measuring properties of celestial bodies or populations for which some distinctive properties are already available. Consider, for example, a subpopulation of galaxies found to exhibit Seyfert-like spectra in the optical band (property A) that have already been examined for nonthermal lobes in the radio band (property B). Then the conditional probability that a galaxy has a Seyfert nucleus given that it also has radio lobes is given by equation (2.11), and this probability can be estimated from careful study of galaxy samples. Similar inferences can be made with

Bayes' Theorem

Let B_1, \dots, B_k be a partition of the sample space Ω .

If A is any event in Ω , then to compute

$P(A)$, one can use probabilities of pieces of A on each of the sets B_i and add them together to obtain the Law of Total Probability

$$P(A) = P(A|B_1)P(B_1) + \dots + P(A|B_k)P(B_k). \quad (2.13)$$

For B_k possible outcomes, using the definition of conditional probability and the Law of Total Probability,

$$P(B_i | A) = \frac{P(A | B_i)P(B_i)}{P(A | B_1)P(B_1) + \dots + P(A | B_k)P(B_k)}.$$

Bayes' Theorem is thus just a necessary result of probability theory, or logic, based on the three Axioms. In Bayesian inference the terms are given a specific meaning (MSMA, p.63-64)

$$P(M_i(\theta) | X) = \frac{P(X | M_i(\theta))P(M_i(\theta))}{P(X | M_1(\theta))P(M_1(\theta)) + \dots + P(X | M_k(\theta))P(M_k(\theta))}$$

Let X represent the data, and M represent the space of models or hypotheses that depend on parameters theta

$P(M_i | X)$ = conditional probability of M_i given X = posterior probability density

$P(X | M_i)$ = conditional probability of X given M_i = likelihood function

$P(M_i)$ = prior or marginal probability of M_i = prior information

$P(X)$ [denominator] = marginal probability of X = normalizing constant

Bayes' Theorem in English:

The posterior distribution of a chosen model given the data is equal to the normalized product of (the likelihood of the data for that model) and (the prior probability that the model is true without reference to the data)

Posterior distribution

Once the prior distribution and alternative hypotheses (range of θ parameters) are specified, and the data are obtained, the posterior distribution can be calculated. This distribution can be plotted in the p -space of its parameters. These plots give information on any non-Gaussianity and multimodality of the posterior. Typically, the scientist is interested in the 'best' Bayesian estimator for the parameters; i.e. the maximum (mode) of the posterior. The *credible region* around this value is then estimated.

Prior distributions

This is often the controversial aspect of Bayesian inference, because subjective judgment or simplistic uninformative priors are often used. For uniform priors, maximizing the posterior often gives the same result as maximum likelihood estimation, although interpretation of results differ. Bayesian inference is most effective when the scientist **wants** to bias the likelihood based on the data using scientifically meaningful prior constraints on the the parameters.

A simple astronomical example

Is this new active galactic nucleus radio-loud?

Let X be a random variable taking two values: $X=1$ indicates *Yes* and $X=0$ is *No*.
Let θ be a parameter denoting AGN radio-loudness: θ_1 indicates *Yes*, θ_2 is *No*

From previous AGN surveys, the astronomer expects a probability of radio-loudness: $P(\theta = \theta_1) = 0.1$.

The new AGN under study was observed with a radio telescope sensitive enough to measure radio-loudness 80% of the time in radio-loud AGN:
 $P(X=1 | \theta_1) = 0.8$. However, 30% of the time the telescope detects irrelevant radio emission from star formation in the host galaxy: $P(X=1 | \theta_2) = 0.3$.

Use Bayes' Theorem to calculate the chances than an AGN with detected radio emission is truly a radio-loud AGN:

$$P(\theta = \theta_1 | X = 1) = \frac{\text{likelihood} \times \text{prior prob}}{\text{sum of marginal probs}} = \frac{0.8 \times 0.1}{0.8 \times 0.1 + 0.3 \times 0.9} = \frac{0.08}{0.35} = 0.23.$$

Wow! Only 23% of true radio-loud AGN are clearly identified in this survey: 77% are either false negatives or false positives. Trying different assumptions shows that the result is moderately sensitive to the value of the prior (0.10) but is very sensitive to the false positive fraction (0.30). If this is reduced to 0.05, then the discovery fraction of true radio-loud AGN rises to 95%. Bayesian calculations can help the astronomer evaluate how the science goals can be better achieved in a future experiment.

Prior distributions: Uninformative priors

Many Bayesian studies (in astronomy and elsewhere) do not have an empirical or subjective basis for specifying the distribution of a model parameter in advance of the experiment/observation at hand. In such cases, an uninformative prior is used to weight the likelihood in Bayes' Theorem. There are two classes commonly used.

Conjugate priors In the 20th century before MCMC methods, statisticians often chose posteriors and priors from conjugate families, as this greatly aided analytic calculations. Example:

Gaussian model → normal inverse Gamma distribution prior

model $Y_i \sim N(\mu, \sigma^2), i=1, 2, \dots, n$

likelihood $L(\mu, \sigma^2) = k \sigma^{-n} \exp(-\sum (Y_i - \mu)^2 / 2\sigma^2)$

conjugate priors $\sigma^2 \sim \Gamma(c, d^{-1}) \quad \mu \sim N(a, b^{-1}, \sigma^2)$

Binomial model → beta distribution prior

Poisson model → Gamma distribution prior

Pareto model → Gamma distribution prior

etc. (see http://en.wikipedia.org/wiki/Conjugate_prior)

Uninformative priors:

These priors make few or no assumptions about the distribution of model parameters. Two common choices:

- The uniform distribution over the full space of possible values. This often reproduces results from maximum likelihood estimation.
- Jeffreys prior $\pi(\theta) = |I(\theta)|^{1/2}$, I is the Fisher Information Matrix

However, use of uninformative priors is controversial and many statisticians do not support their use:

1. Many are **improper priors** that do not integrate to unity (often the integral is infinite). Thus they are not p.d.f.'s.
2. The results depend on arbitrary choices. In an astrophysical model, is the prior of X or $\log(X)$ assumed to be uniform? For the normal model, is the prior of the variance or the standard deviation assumed to be uniform? A uniform s.d. allows Bayesian calculations to reproduce many classical results. These are based on mathematical convenience, and do not really represent any PRIOR knowledge. ['Arbitrary' here might be phrased 'garbage in -- garbage out']

Proper use of priors

Many astronomers proceed with a uniform or other uninformative improper prior without consideration of alternatives. We discourage this practice. A reasonable alternative is to try different reasonable proper priors and, if the results are compatible, report them as scientifically reliable results.

When flat or uninformative priors are used together with estimation using the mode of the posterior (MAP or HPD best fit), then we recommend that the Bayesian approach be dropped and the Maximum Likelihood Estimation formulation be used instead.

When the prior can be reliably established from detailed scientific information available from earlier observations or from astrophysical theory, then we encourage use of these informative priors with a Bayesian approach. However, it is wise to examine the relative influence of the prior and the data on the scientific result for the particular situation at hand.

A published Bayesian analysis should communicate the prior distribution in sufficient detail that other scientists can apply it to their data.

A worried viewpoint about uninformative priors

“Because the prior is inescapably part of the model in the Bayesian approach, marginal likelihoods, Bayes factors and posterior model probabilities are inescapably sensitive to the choice of prior. In consequence, it is only when those priors that differ between alternative models are really precise and meaningful representations of prior knowledge that we can justify using Bayes factors and posterior model probabilities for model selection. Even then the computation of the marginal likelihood is often difficult.”

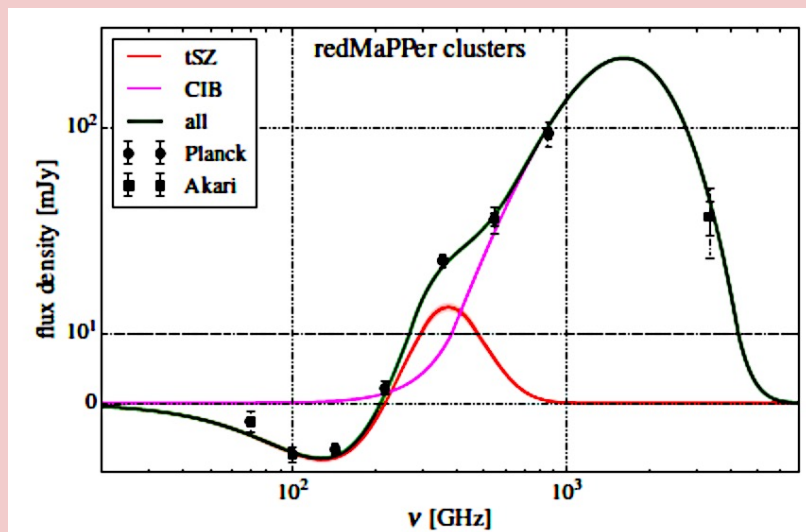
Simon Wood, *Core Statistics* (2015)

Bayesian posterior

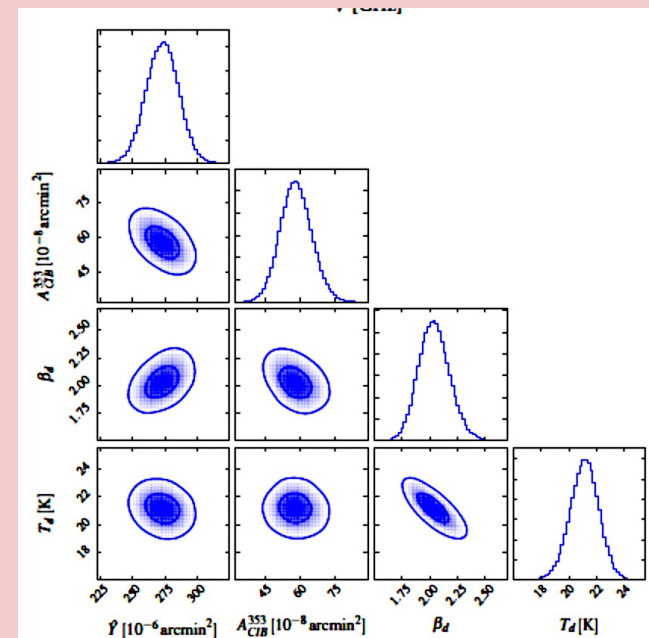
The result of a calculation of Bayes' Theorem for a dataset and a model space is the distribution of the posterior. Astronomers often plot univariate and bivariate projections of a multivariate posterior estimated by MCMC sampling.

Example: Model of Sunyaev-Zel'dovich distortion to the cosmic microwave background spectrum (taking the dust-induced cosmic infrared background (CIB) variations into account) applied to 26,111 galaxy clusters from the Sloan Digital Sky Survey (Soergel et al. 2017)

Best fit model

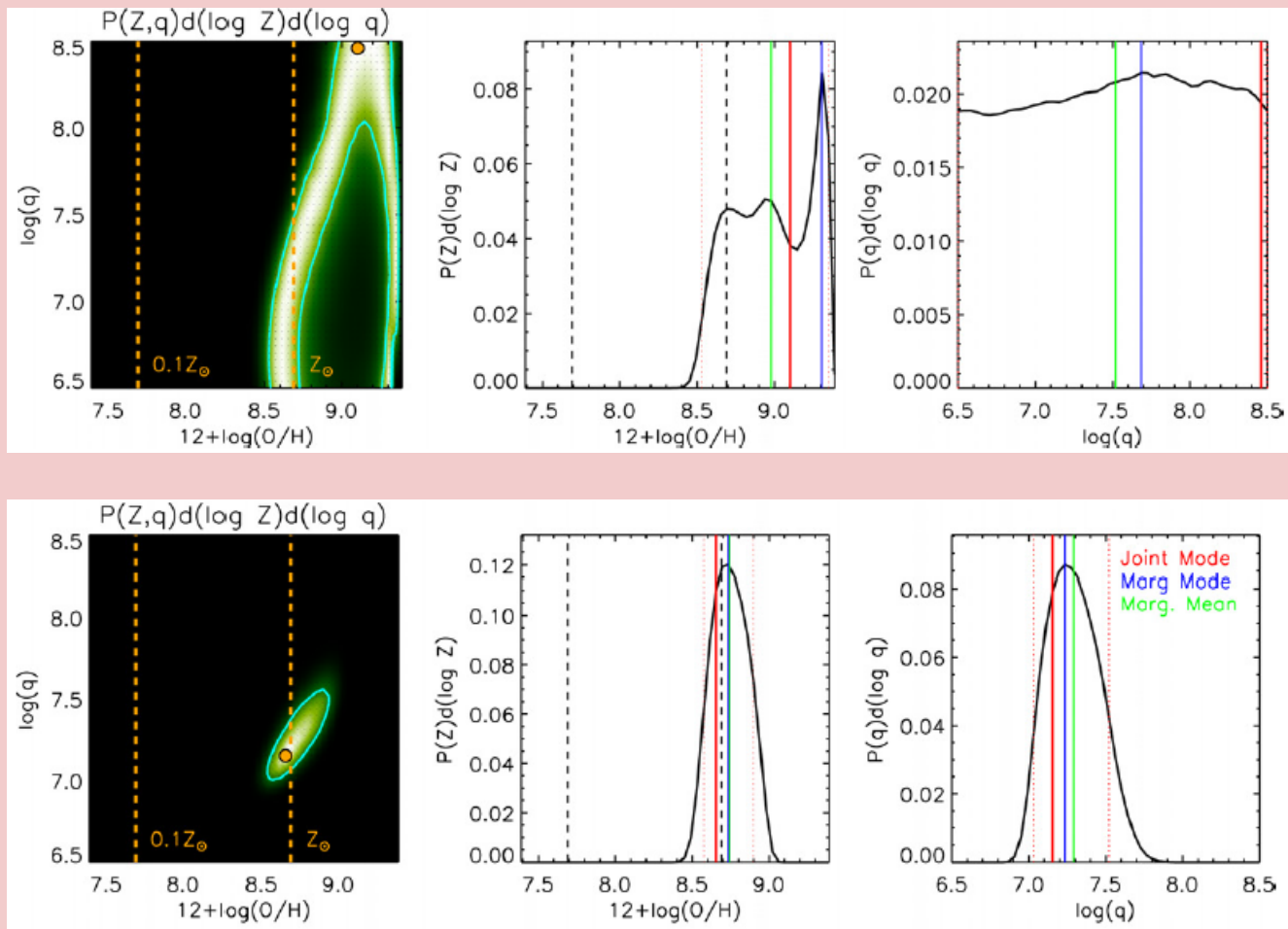


Marginal posterior distributions



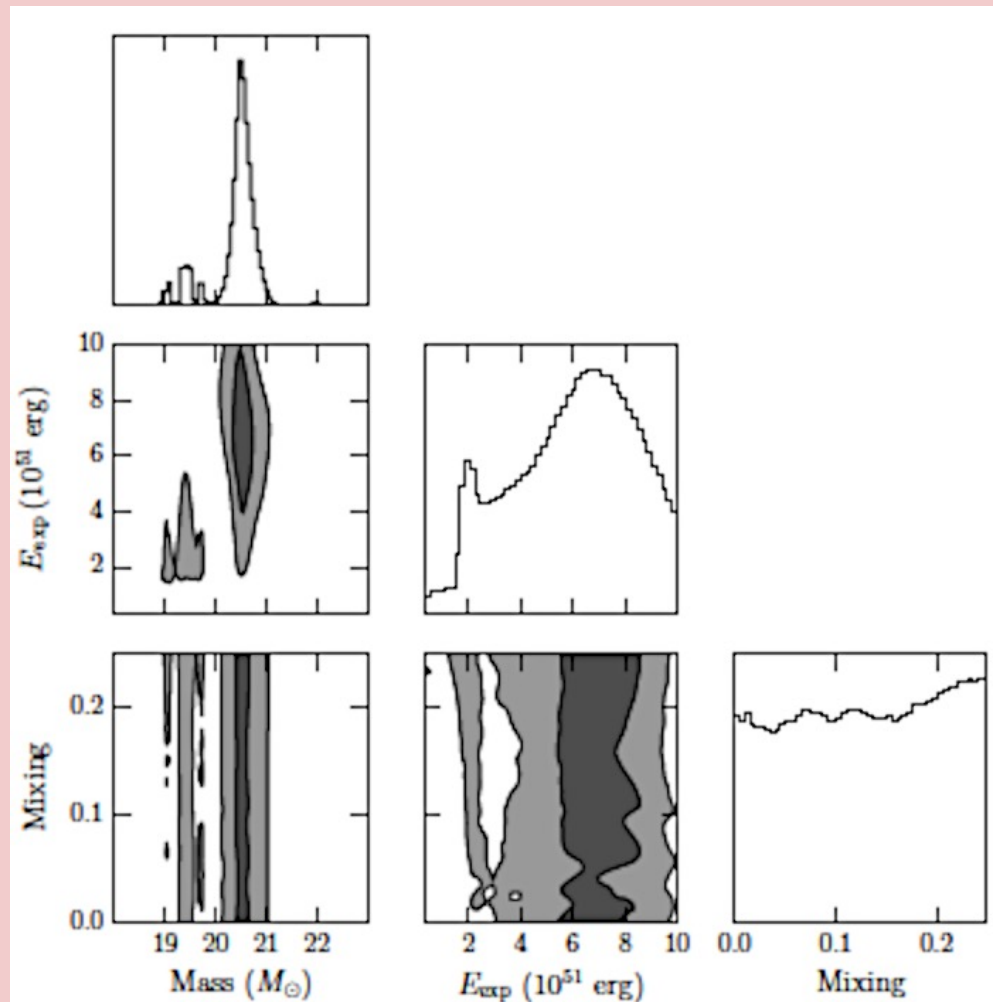
Some messier posteriors

Bayesian analysis of ionization and metallicity in HII regions. The top panels show joint and marginal posteriors of ionization and O/H abundance using 2 emission lines. The bottom panels show the posterior using 8 emission lines.



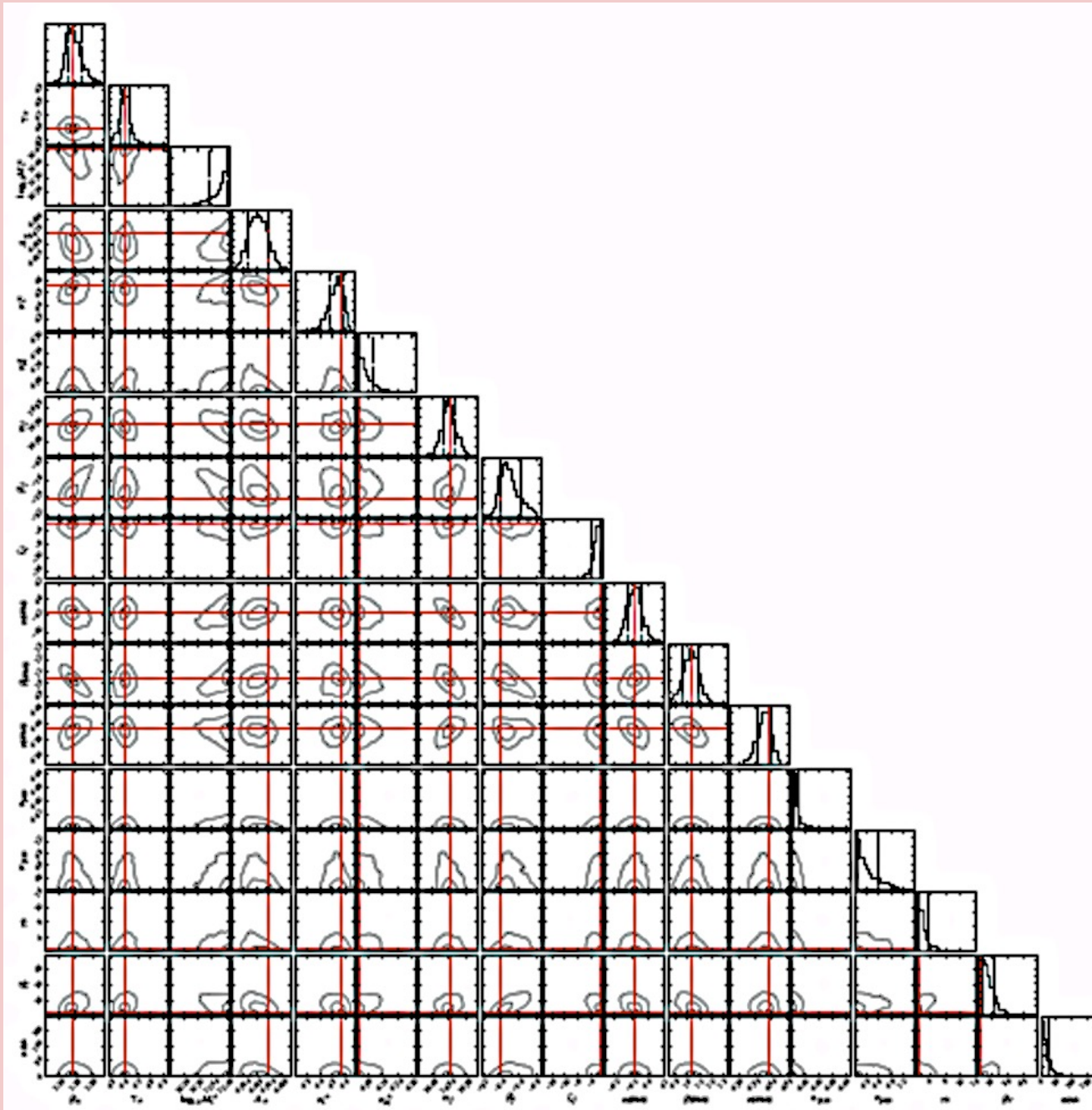
(Blanc et al. 2015)

Modeling a high-redshift metal-poor damped Ly α absorption (DLA) system



Cooke et al. 2017

A more complicated posterior



Bayesian parameter estimation

(or Summarizing the posterior distribution)

Scientists often seek a single 'best' model giving 'best-fit' parameters for the dataset and the model space, rather than a multivariate distribution of model probabilities.

Three approaches are commonly used in Bayesian inference. The choice should be based on the a previously specified 'loss function' (or 'risk function') that quantifies the scientific value of alternative models. The principles arise from Bayesian decision theory, a branch of information theory.

- The **mode** of the posterior distribution. This is sometimes called the **maximum a posteriori (MAP)** or the **highest posterior density (HPD)** estimate. For uniform priors or very large datasets, the posterior mode gives model parameter values approaching the classical maximum likelihood estimators. For an informative prior, the MAP solution is a weighted average of the MLE of the prior and the MLE of the data. In decision theory, the mode is preferred when the cost of a wrong answer is high (posterior loss is binary).
- The **median** of the posterior is preferred when the cost of a wrong answer is low (posterior loss scales as the linear 'distance' between models)

- 3. The **mean** of the posterior distribution. This a weighted mean of the likelihood and the prior:

$$E(\theta|\mathbf{X}) = \frac{\int_{\Omega} \theta L_{\mathbf{X}}(\theta) \pi(\theta) d\theta}{\int_{\Omega} L_{\mathbf{X}}(\theta) \pi(\theta) d\theta}$$

This is simply the expectation of the posterior distribution. The mean is preferred for intermediate cost functions (posterior loss scales as the squared distance between models).

J. O. Berger, *Statistical Decision Theory and Bayesian Analysis*, 2nd ed 1985

C. Robert, *The Bayesian Choice: From Decision-Theoretic Foundations to Computational Implementation*, 2nd ed 2007

Unfortunately, astronomical research does not usually have a clear loss function, so astronomers are subjectively choosing to report medians, means and modes.

Essentially, the research community must choose a consistent summary statistic, much as it chooses a consistent significance level (3-sigma) for reporting results of hypothesis tests.

*Some experts have suggest the **posterior median** for general use. Other experts suggest avoiding summarizing the posterior and use/discuss the entire distribution.*

Bayesian credible intervals

The Bayesian **credible interval** of parameter values (or **credible region** for p -dimensional models) around the MAP value can be estimated from the analytical or computed posterior distribution. This plays the role of the *confidence interval* in classical statistical inference. The credible interval can be found by solving for lower and upper functions such that

$$P(L(\mathbf{X}) \leq \theta \leq U(\mathbf{X}) \mid \mathbf{X}) = 1-\alpha$$

where $\alpha=0.05, 0.01$, etc. For simple cases (e.g. Gaussian mean with a conjugate prior), the credible interval can be computed analytically. In realistic cases, it is computed numerically by computing values of the posterior distribution around the best fit value.

Bayesian model selection

An important class of hypothesis tests is **model selection**, the comparison of two alternative models for a given dataset. When applied to nested models, this problem is important for deciding how many parameters is needed to adequately fit the data in a parsimonious fashion.

Bayesian model selection is based on the **Bayes factor**, or ratio of posteriors, given by

$$B_{12} = \frac{P(\mathbf{X}|M_1, \pi)}{P(\mathbf{X}|M_2, \pi)}$$

The ratio of the probabilities of the two models, or **odds ratio**, is the product of the ratio of likelihoods and the Bayes factor:

$$O_{12} = \frac{P(M_1|\pi)}{P(M_2|\pi)} B_{12}$$

Note that the Bayesian odds ratio is equal to the classical likelihood ratio test (LRT) when the priors for the two models are equal. This is often the case when the priors are uninformative (i.e. there is little prior knowledge about the model parameters).

History: The LRT was established by theorem by Neyman & Pearson (1933), and Wilks (1938) showed it asymptotically follows a chi-squared distribution.

Model selection is an example of Bayesian hypothesis testing and has a number of advantages over classical hypothesis testing:

- The Bayes factor automatically accounts for the number of parameters, favoring the more complicated model only if the ratio of the likelihoods is sufficiently high. In classical MLE, the penalty for model complexity is debated, and does not arise naturally from the mathematics.
- Both classical and Bayesian analysis often use the Bayesian Information Criterion (BIC), which is an approximation to the Bayes factor.
- Bayes factors allow comparison of nonnested models, and Bayesian model averaging can be used to account for model uncertainty.

To compute the Bayes factor, we need to calculate the *marginal likelihood* of each model for the available data

$$f(\mathbf{y}|M_1) = \int f(\mathbf{y}|\boldsymbol{\theta}_1)f(\boldsymbol{\theta}_1)d\boldsymbol{\theta}_1$$

Use of these marginal likelihoods in model selection accounts for differences in model complexity – models with larger ‘volumes’ of parameter space are not automatically favored, as they are when the LRT is used.

However, two difficulties arise. First, it may be hard to compute the marginal likelihoods for all parameters and (possibly) for a wide range of models. Second, the marginal likelihoods are sensitive to the prior and will change values for different uninformative priors unless the same improper priors are used in all models.

In practice, it is often easier to compute the approximate than the full odds ratio.

Bayesian marginalization

Many problems have variables or parameters of little scientific interest (e.g. detector background vs. astronomical signal). Bayesian formulations allow direct **marginalization** (integration) over nuisance variables.

Consider a case of a vector of k parameters where we are interested in the i^{th} parameter and not the others:

$$P(\theta_i) = \int d^{k-1}\theta P(\theta_1, \theta_2, \dots, \theta_k | x_1, x_2, \dots, x_n)$$

The distribution of the interesting parameter θ_i takes into account information about all of the other parameters. This 'averages out' the influence of other parameters.

Hierarchical Bayes' modeling

When a Bayesian model is based on a prior distribution that itself has unknown parameters, the calculation must simultaneously solve for the model parameters and the prior **hyperparameters**. This is called a **hierarchical Bayes' model**. The acquisition of additional data simultaneously updates the prior distribution and constrains the model parameters of interest. Examples:

- the prior is a mixture of two distributions with the mixing fraction serving as a hyperparameter
- the model parameters θ are generated from a process governed by a hyperparameter ψ . Then (ignoring normalizations)

$$P(\theta, \phi | \mathbf{X}) = P(\mathbf{X} | \theta) P(\theta, \phi) P(\phi)$$

Hierarchical Bayesian models are increasingly common in the astronomical literature.

Philosophical considerations for model fitting: Comparing classical and Bayesian approaches

Classical statistical procedures, e.g. as developed by R.A. Fisher, consider probabilities in the context of real or hypothetical random experiments where a measurement of some property or sample is made many times, each time with some error. A Fisherian 'frequentist' accepts the model parameters which maximizes the likelihood that the data satisfy the model.

Bayesians assign probabilities to many situations, e.g. subjective decisions as well as models based on physical measurements. Here the randomness is associated not with the measurement, but with the prior beliefs regarding the question under study and with the models used to interpret the findings.

A Bayesian views probability as the plausibility of a situation or interpretation based on a combination of current and past information.

A frequentist views probability as the chance of a situation or interpretation assuming many hypothetical experiments were made, without consideration of past information.

Bayesian calculations average over model space, while frequentist calculations averages over sample space (and optimize over model space):

Bayesian: Data are fixed, hypotheses vary

Frequentist: Hypothesis is fixed, data vary

'Why isn't every physicist a Bayesian?' (particle physicist R. Cousins)

- For many simple situations, frequentist and Bayesian solutions are (nearly) identical
- Frequentist estimation is typically simpler mathematically and computationally. Except for trivial problems, Bayesian estimation requires arduous computation for the calculation of posteriors in the full space of possible parameter values
- Bayesian estimation will be biased if the prior distributions are misspecified. If priors are not known, MLEs may be preferred.

- Bayesian estimation uses and are more informative when prior (i.e. not involving the data at hand) information about the model is available. This prior information can arise from astrophysical theory and/or previous empirical study.
- Bayesian model selection has advantages over frequentist model selection. The Bayes Factor and BIC can evaluate evidence in *favor of* (not just against) a model; be applied to non-nested model alternatives; incorporates external (prior) information; and has a natural compensation for model complexity (Occam's Razor). However, Bayesian model selection does not give formal probabilities.
- MCMC is not intrinsically related to Bayesian inference. They can (and should) be used in MLE to map and characterize (unweighted) likelihood surfaces.

Some disadvantages of Bayesian inference

Bayesian inference depends on the specification of a large, and often ill-defined, space of possible models.

For each possible model, Bayesian inference requires quantitative statement of the distribution of each models of interest prior to the acquisition of data. This often gives a subjective element to the procedure.

Bayesian model fitting requires specification of, and integration over, a universe of alternative theories. This is often both conceptually and computationally difficult. Simulations may take millions of iterations and may not converge. MLE is much less computationally demanding.

Some advantages of Bayesian inference

When the scientist indeed have prior knowledge (from previous observations or from astrophysical theory) of the parameter distributions, this can readily be incorporated into the Bayesian prior. Bayesian inference takes full account of this information. Such ancillary information and can only be included into frequentist calculations with difficulty.

Bayesian 'marginalization' can treat the effects of nuisance variables (e.g systematic error, unobservable or uninteresting variables) with greater transparency than is often achieved with frequentist calculations.

Bayesian hypothesis tests treat hypotheses symmetrically, and Bayesian model selection can give probabilities that a model is correct. Classical hypothesis tests only give probabilities that a model is incorrect.

Broad advice on choosing inferential approach

When little is known about a scientific problem and the questions addressed are straightforward, then nonparametric statistics and hypothesis tests may be most appropriate. (MSMA Chpts 3.5 & 5)

When a parametric model, either heuristic or astrophysical, can be reasonably applied, and the experimental situation is relatively simple, then frequentist point estimation may be valuable (least squares & MLE). (MSMA Chpts 3.4 & 4)

Bayesian inference is best pursued when the situation is known to have external constraints (informative priors based on real knowledge from previous astronomical observations or from astrophysical theory), nuisance variables are present, hierarchical relationships exist between variables. (MSMA Chpt 3.8)

Don't hesitate to pursue multiple avenues of analysis

1. Nonparametrics – ‘Let the data speak for themselves’ (Fisher?)
1. Maximum likelihood estimation – Can the data, considered in isolation, be well-fit by a chosen mathematical model? What are the best-fit parameter values? Is the best fit a good fit?
1. Bayesian with simple priors – What influence does prior knowledge about the parameters have on the best-fit solution?
1. Hierarchical Bayes – What can we learn about more complicated models with latent variables, prior hyperparameters, several modeling stages, etc.