



OBSERVATIONAL COSMOLOGY: COSMOLOGICAL INFERENCE

BAYESIAN AND FREQUENTIST STATISTICS

- Probability as frequency:

The classical approach to statistics defines the probability of an event as “the number of times the event occurs over the total number of trials, in the limit of an infinite series of equiprobable repetitions.”

- Probability as degree of belief:

The Bayesian viewpoint is based on the simple and intuitive tenet that: “probability is a measure of the degree of belief about a proposition”.

BAYESIAN AND FREQUENTIST STATISTICS

- **Bayesian:** data are fixed, model is repeatable
- **Frequentist:** model is fixed, data are repeatable

Say $H_0 = (72 \pm 8)$ km/s/Mpc. Then:

Bayesian: the posterior distribution for H_0 has 68% of its integral between 64 and 80 km/s/Mpc. The posterior can be used as a prior on a new application of Bayes' theorem.

Frequentist: Performing the same procedure will cover the real value of H_0 within the limits 68% of the time. But how do I repeat the same procedure (generate a new H_0) if I only have one Universe?

Good references:

Bayesian: R. Trotta, “Bayes in the Sky”, <https://arxiv.org/abs/0803.4089>

Frequentist: Feldman & Cousins, “A Unified Approach to the Classical Statistical Analysis of Small Signals”, <https://arxiv.org/abs/physics/9711021>

Example of one cosmology inference done both Bayesian and frequentist way: G. Efstathiou, “The Statistical Significance of the Low CMB Multipoles”, <https://arxiv.org/abs/astro-ph/0306431>

Or “Robust constraints on tensor perturbations from cosmological data: a comparative analysis from Bayesian and frequentist perspectives” <https://arxiv.org/pdf/2405.04455>

- Bayesian:
 - can give probabilities for models
 - depends on both prior and likelihood (of data)
 - currently the dominant method in cosmology
- Frequentist:
 - doesn't give probabilities of models, only of hypotheses
 - doesn't depend on prior, just likelihood
 - currently the dominant method in particle physics

BAYES' THEOREM

Posterior

What you know
after the
experiment

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

Evidence

Normalization constant

$$\int dp P(d|p, M) P(p|M)$$

Likelihood

What you learn from
the experiment

$$\propto P(d|pM)P(p|M)$$

Prior

What you knew
before the
experiment

Observed data

Parameters

Model

Note: $P(A|B)$ reads “the probability of A given B”

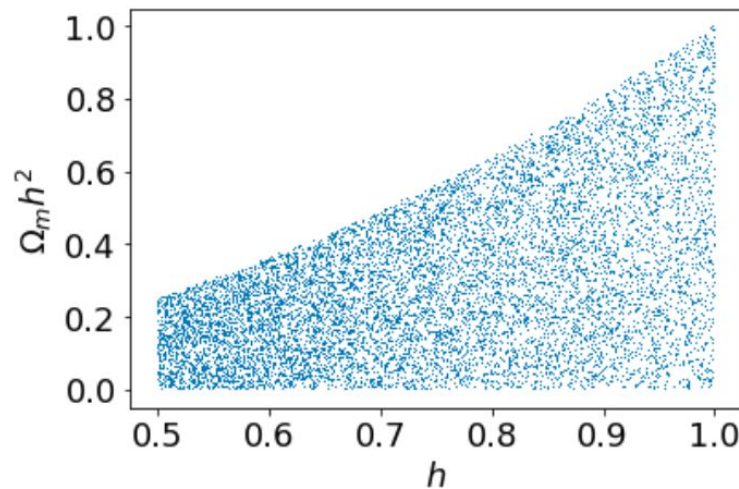
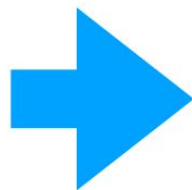
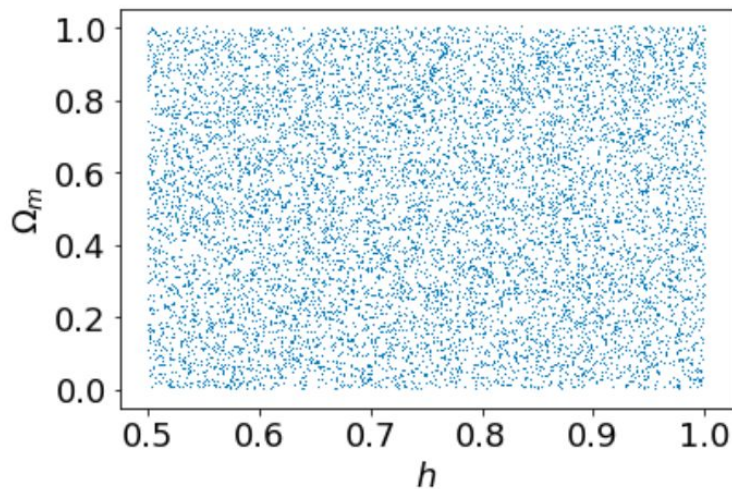
PRIORS

- Priors quantify what you knew about the parameters before you start
 - Theoretical limits, preferences, things that must be true from simpler data
- In regions where your likelihood is zero your prior doesn't matter for parameter estimation, but can for more advanced *model selection*
- It is common practice in cosmology to use uniform priors for most parameters

The rationale is that we should assign equal probability to equal states of knowledge. However, flat priors are not always as harmless as they appear. One reason is that a flat prior on a parameter θ does not correspond to a flat prior on a non-linear function of that parameter, $\psi(\theta)$. The two priors are related by:

$$p(\psi) = p(\theta) \left| \frac{d\theta}{d\psi} \right|$$

PRIORS



Jointly uniform priors on Ω_m - h

Implied priors on $\Omega_m h^2$ - h

E.g. The parameter actually
constrained by CMB data

POSTERIOR ESTIMATION

The challenge: map out a posterior in multi-dimensional parameter space.

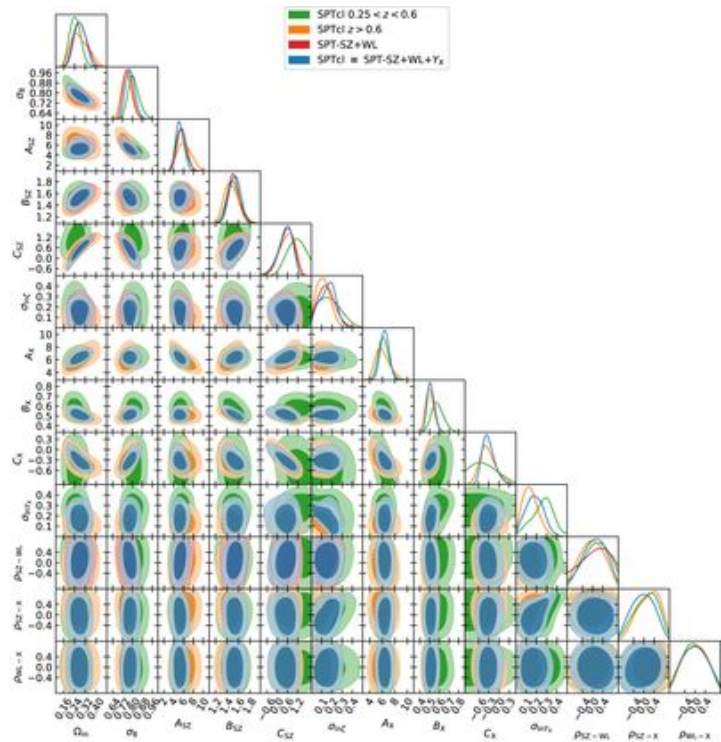
Example: say there are just 10 parameters.
Let's say calculation takes just 1 second/model.
Say you want a grid with 20 values in each par.

Then

$$N = 20^{10} \approx 10^{13}$$

⇒ it would take 300,000 years to do it!

⇒ Totally impossible, ever!!

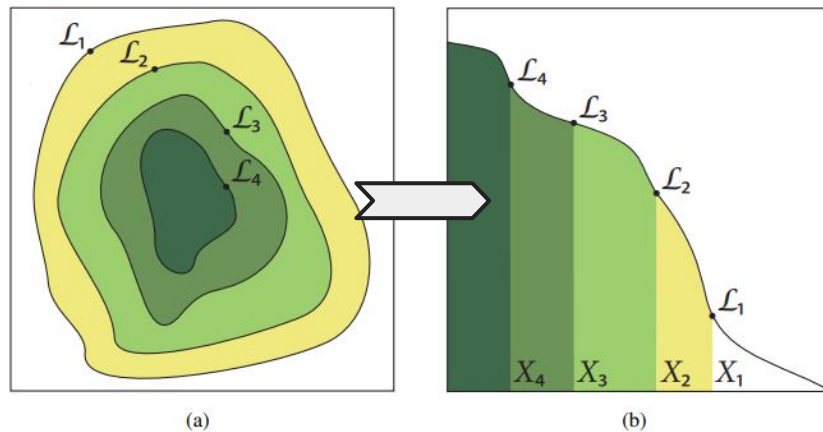
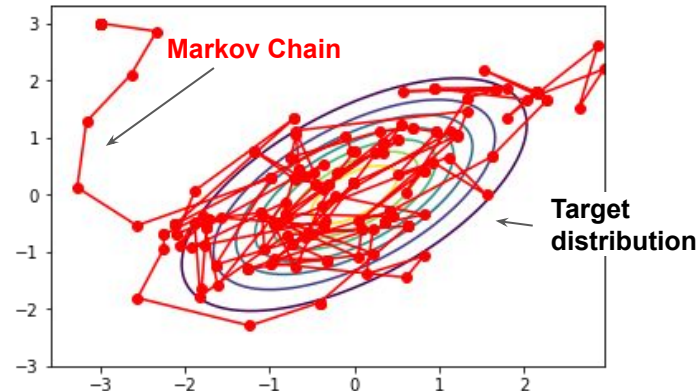


POSTERIOR ESTIMATION

Over the past years many sampling techniques have been developed to overcome this issue (See [this](#) for a review). The general idea is to sample the parameter space in a clever way in order to map out the high-probability volumes. The methods can be divided in:

- Monte Carlo Markov Chains methods: e.g. Metropolis-Hastings (Metropolis+1953), Emcee (Foreman-Mackey+2010)
- Nested sampling methods: e.g. Multinest (Feroz+2009,2013), Polychord (Handley+2015)

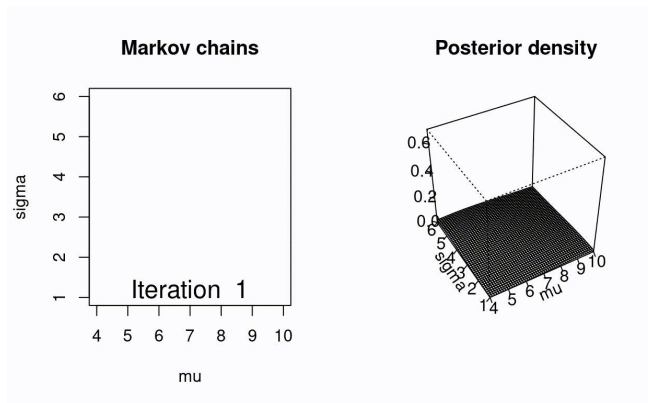
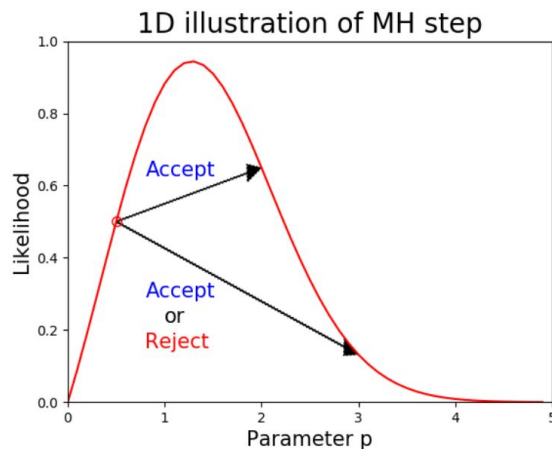
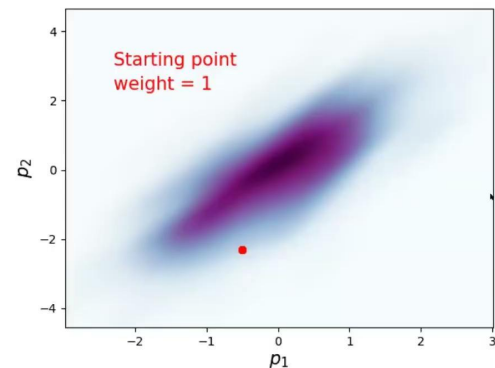
In both case the density of the sampled points is proportional to the parameter posterior we seek to estimate



MCMC: METROPOLIS-HASTING ALGORITHM

- ▶ at step t , at some parameters p_t
- ▶ propose move to $p_t' = p_t + \Delta p_t$ (randomly draw Δp_t)
- ▶ evaluate $r = L(p_t')/L(p_t)$
- ▶ MH step:
 - ▶ if $r > 1$ **accept move**
 - ▶ if $r < 1$ generate a random number $\alpha \in [0, 1]$
 - ▶ if $\alpha < r$, **accept move**
 - ▶ if $\alpha > r$, **reject move**
- ▶ $t=t+1$

One can prove that,
with this rule,
one asymptotically recovers the
true posterior



NESTED SAMPLING METHODS

CORE IDEA: Instead of sampling the entire prior space (like MCMC), NS sequentially shrinks the sampling region by discarding low-likelihood points, concentrating on high-probability regions.

Algorithm:

1. Sample N point from the prior distribution
2. At each step i -th:
 - a. Find the point with the lowest likelihood L_i
 - b. Replace it with a new point having $L > L_i$
 - c. Record the “dead point” and its likelihood
3. Stop when the remaining live point contribute negligibly to the evidence

Output:

1. Posterior samples
2. Evidence estimate

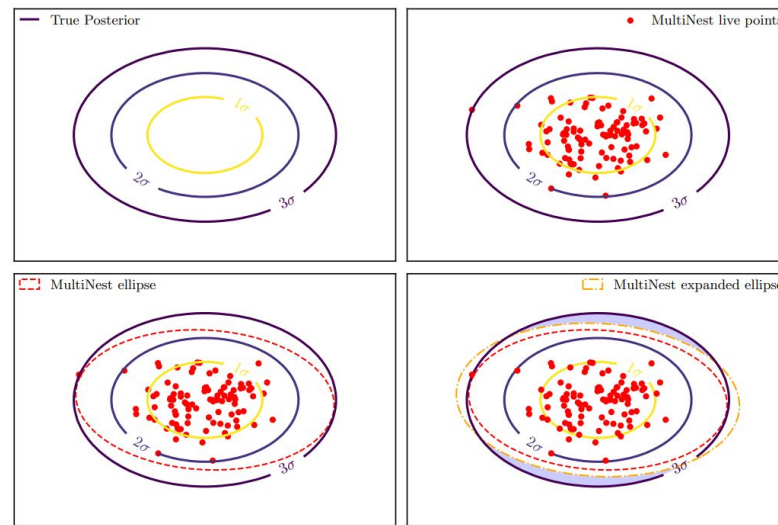
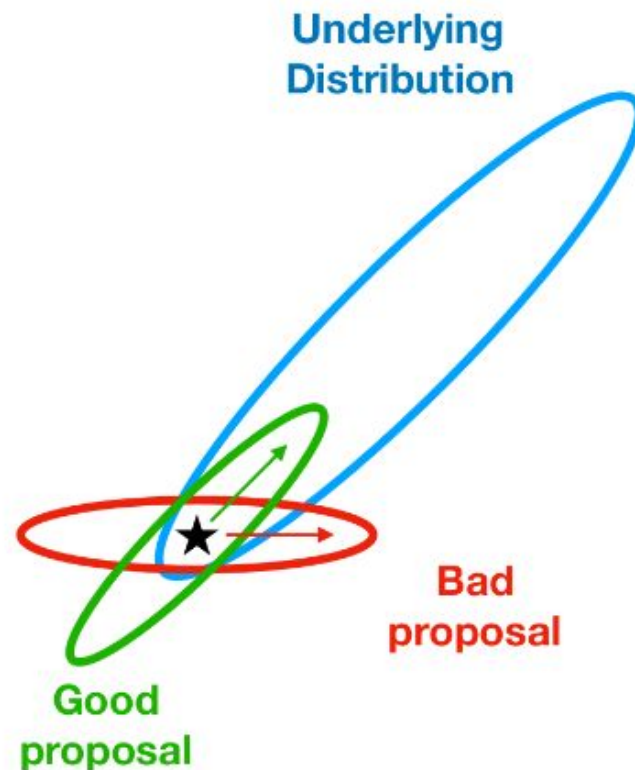


Figure 1. An example of MultiNest’s ellipsoidal sampling, and how it can lead to biases. When trying to sample a certain distribution (top left), MultiNest randomly generates some points (top right). It then uses the covariance matrix obtained from those points to calculate an ellipsoid enclosing all existing live points (bottom left, dashed line). That ellipsoid is expanded in volume by a factor inversely proportional to the efficiency, and samples are drawn from that ellipsoid (bottom right, dot-dashed line). As the latter plot shows in the light blue regions, if the magnification factor is not big enough (i.e. the efficiency is too high), this can lead to a bias in the estimation of the evidence.

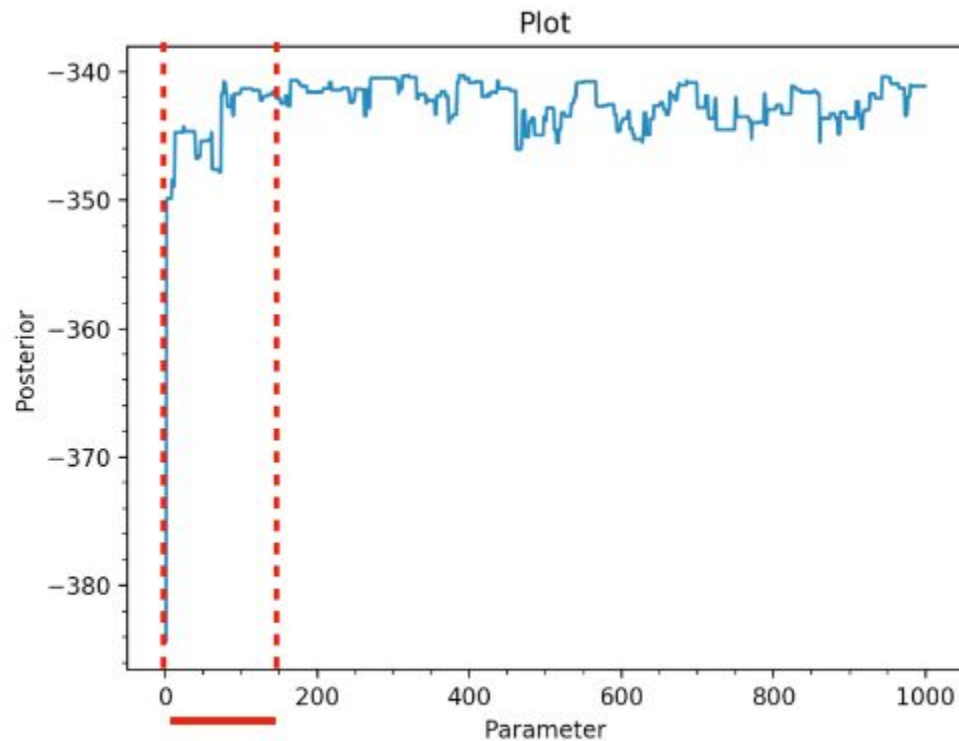
SAMPLING THE PARAMETER SPACE

- Efficiency of MH depends dramatically on how good the proposal is
- A bad proposal will not converge in any practical length of time
- The ideal proposal matches the shape of the underlying distribution
 - We don't know this, but can look for best approximation



BURN-IN

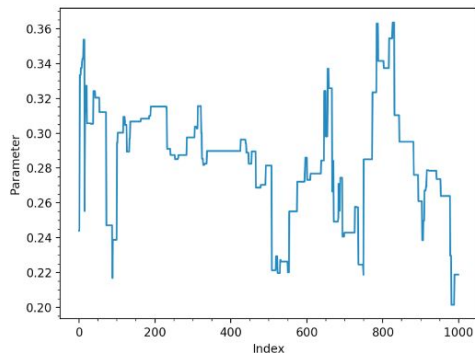
- Unless you're doing a simulation where you know the truth, unlikely to start at the best-fit value
- Will take some iterations to get near this point
- Need to exclude these



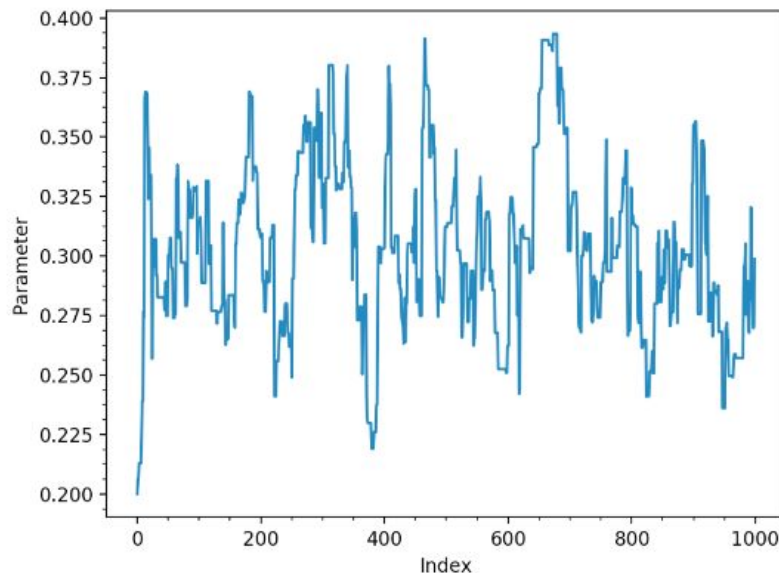
Burn-in - exclude from sampling

CHECKING CONVERGENCE

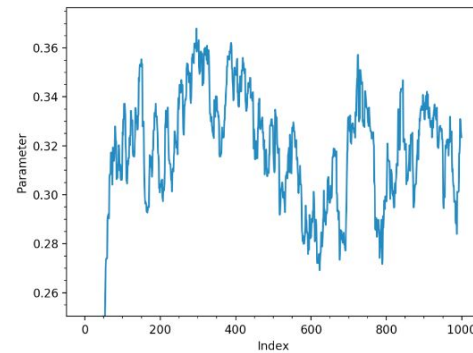
- Good MH chains look like white noise if you plot one parameters values throughout the chain



Bad - not long enough.
Chain getting stuck for long periods
suggests covariance too large.
Acceptance rate too low.



Looks reasonable - could be a bit longer



Bad - not long enough.
Chain is random walking, taking long divergences
from mean suggests covariance too small.
Acceptance rate too high.

INTERPRETING THE OUTPUT

WEIGHT	P_1	P_2	P_3	...	P_N
5	0.2	-0.3	0.15	...	2.8
1	-0.7	0.4	0.12	...	3.5
12	0.7	0.1	0.19	...	1.7
...
...

(~ MILLION ROWS)

To get the posterior probability,
simply histogram the parameter values vs weights - this is your posterior!

Want to look at posterior in p_3 marginalized over all other parameters?
Simply plot histogram of p_3 values vs weight (easy!)

GOODNESS OF FIT

The goodness of fit is often estimated from the best-fit parameter values using a χ^2 statistic (which is formally correct only for Gaussian distributions):

$$\chi_{\text{Best-fit}}^2 = (\vec{d} - \vec{m}(\theta_{\text{BF}}))^T C^{-1} (\vec{d} - \vec{m}(\theta_{\text{BF}}))$$

where C is the data covariance matrix. *This method does not account for the uncertainties on the estimated parameters ϑ .*

To assess the goodness of fit from the $\chi_{\text{best-fit}}^2$ one computes $p(\chi^2 > \chi_{\text{BF}}^2 | \nu)$ the probability to exceed the χ_{BF}^2 , assuming a χ^2 -distribution with ν degree of freedom: $\nu = N$. Data points - N. effective parameters. The number of effective parameters, for correlated parameters and/or for a prior-informed analysis, is smaller than the total number of free parameters.

Interpreting Reduced Chi-Square

- ($\chi_{\text{red}}^2 \approx 1$): The model is a good fit for the data. (Residuals are of the same order as the uncertainties.)
- ($\chi_{\text{red}}^2 > 1$): The model might not fit well.
- ($\chi_{\text{red}}^2 \gg 1$): Indicates systematic errors, underestimation of uncertainties, or a poor model.
- ($\chi_{\text{red}}^2 < 1$): The model may be overfitting or uncertainties might be overestimated.

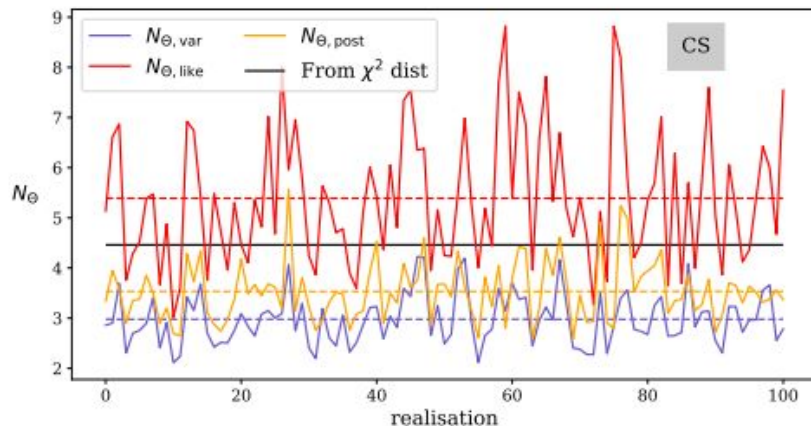
EFFECTIVE NUMBER OF D.O.F.

Effective number of constrained parameters:

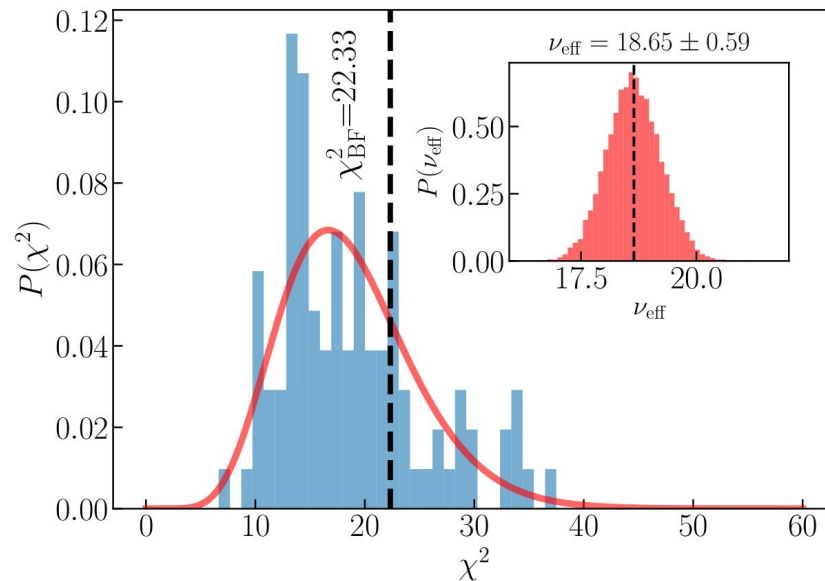
$$N_{\Theta, \text{var}} = 2 \left[\left\langle (\chi^2)^2 \right\rangle_{\text{Pr}(\Theta|d)} - \left\langle \chi^2 \right\rangle_{\text{Pr}(\Theta|d)}^2 \right];$$

$$N_{\Theta, \text{like}} = \left\langle \chi^2 \right\rangle_{\text{Pr}(\Theta|d)} - \chi^2_{\min};$$

$$N_{\Theta, \text{post}} = \left\langle \chi^2 \right\rangle_{\text{Pr}(\Theta|d)} - \chi^2 (\text{argmax} [\text{Pr}(\Theta|d)]),$$



From Joachimi et al 2021



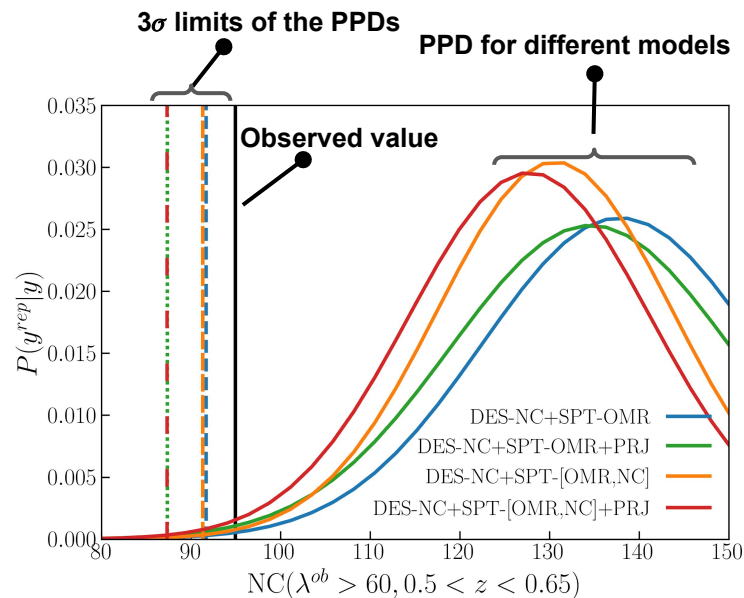
Distribution of the best-fit χ^2 values recovered from 100 mock data realizations generated from the best-fit model of the data. The red histogram in the inset plot shows the posterior distribution for the effective number of degrees of freedom obtained by fitting a χ^2 to the histogram. (DES Collaboration 20)

GOODNESS OF FIT

A more rigorous way to assess the goodness of fit which account for both the data and model uncertainty rely on the **Posterior Predictive Distribution**:

$$P(y^{\text{rep}}|y) = \int d\theta \underbrace{P(y^{\text{rep}}|\theta)}_{\text{Likelihood}} \underbrace{P(\theta|y)}_{\text{Parameter posteriors}}$$

The method consists of drawing simulated values from the posterior predictive distribution of replicated data and comparing these mock samples to the real data to assess their likelihood to be observed (see e.g. Doux+2021)

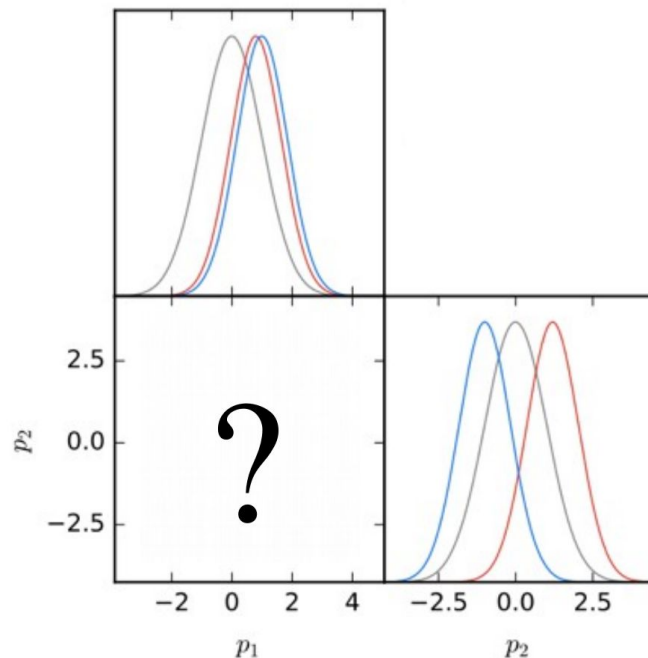


PPD for the observed cluster count in the highest λ/z bin of the DES Y1 data for 4 different model (Costanzi+21)

TENSION BETWEEN DATA SETS

Asses the level of tension (or agreement) between posteriors derived from different data sets might not be trivial in a multi-dimensional parameter space.

E.g. 1d marginalized posteriors which seem to be in agreement ...

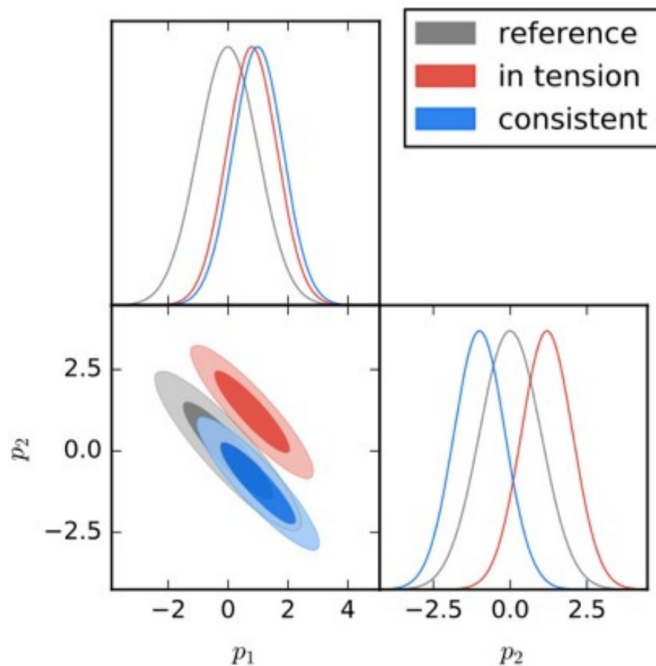


Credit A. Saro

TENSION BETWEEN DATA SETS

Asses the level of tension (or agreement) between posteriors derived from different data sets might not be trivial in a multi-dimensional parameter space.

E.g. 1d marginalized posteriors which seem to be in agreement ...



Credit A. Saro

... might hide tensions in higher dimension space due to “projection effects”

TENSION METRICS

There is no a unique “metric” to assess the level of tension/agreement between data sets, and there exist a number of technique which can be roughly splitted in:

- **Evidence-based methods** seek to answer the question:

Given hypothesis H_1 : ‘The assumed model is capable of generating the data observed by both experiments’, and hypothesis H_2 : ‘The assumed model is not capable of generating the data observed by both experiments’, which hypothesis is preferred by the data under the assumed model’?

- **Parameter-space methods** seek to answer the question:

What is the statistical significance of the differences between the posteriors for experiments A and B, within the parameter space analyzed by both experiments?

(Lemos+2020; see also e.g. Grandis+16, Charnock+17, Raveri+20)

Require the computation of the evidence:

$$P(\vec{d}) = \int d\theta \mathcal{L}(\vec{d}|\vec{\theta}) P(\vec{\theta})$$

In general can be computed directly from the parameter posteriors.
Require good sampling of the tails of the distributions

TENSION METRICS

Bayes Ratio

The Bayes ratio R is defined for independent datasets A and B and for their combination AB as [210]:

$$R \equiv \frac{Z_{AB}}{Z_A Z_B}, \quad (\text{E1})$$

where

$$Z_D \equiv P(\mathbf{D}|\mathcal{M}) = \int d\Theta \mathcal{L}(\mathbf{D}|\Theta, \mathcal{M})\pi(\Theta|\mathcal{M}). \quad (\text{E2})$$

In that expression, z_D is the Bayesian Evidence, L is the likelihood of observing the data given model M and parameter values Θ , and π is the prior probability of those parameters given the model. R can be viewed as a hypothesis test assessing the odds of both datasets being described with a single set of parameters (Z_{AB}) as opposed to two independent sets of parameters ($Z_A Z_B$)

Smaller values of R indicate stronger evidence of tension between measurements from datasets A and B

Jeffrey's scale	
$\ln R < -2.3$	Strong Tension (10:1 odds)
$-2.3 < \ln R < -1.2$	Substantial tension (3:1 odds)
$\ln R > -1.2$	Agreement

Caveat: the value of $\ln R$ depends strongly on the choice of parameter prior ranges

TENSION METRICS

Parameter difference technique:

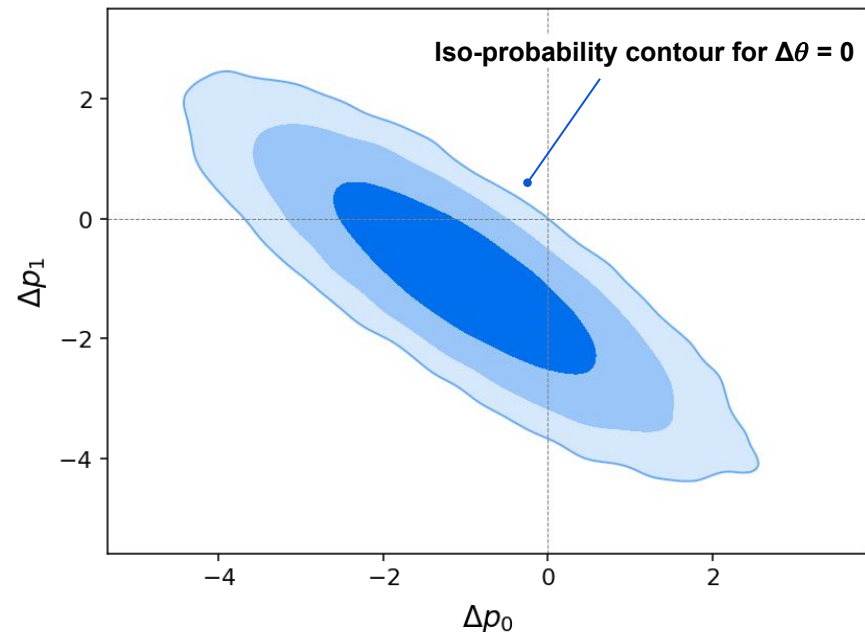
i) Compute the parameter difference probability distribution:

$$\mathcal{P}(\Delta\theta) = \int_{V_p} \mathcal{P}_A(\theta) \mathcal{P}_B(\theta - \Delta\theta) d\theta$$

ii) Determine posterior mass above the iso-probability contour for no shift, $\Delta\theta = 0$

$$\Delta = \int_{\mathcal{P}(\Delta\theta) > \mathcal{P}(0)} \mathcal{P}(\Delta\theta) d\Delta\theta$$

The advantage of this technique is that it can be readily computed directly from the MCMC chains of experiment A and B



Toy model for a two parameter difference distribution. Credit M. Raveri

MODEL SELECTION

To determine which model is preferred by a given data set a simple comparison of χ^2 s might not be sufficient (e.g. if the two models have a different number of parameters, or different priors)

Two widely used techniques for model selection are the evidence ratio and deviance information criterion:

Bayes Evidence Ratio:

$$\frac{P(M_1|\vec{d})}{P(M_2|\vec{d})} = \frac{P(\vec{d}|M_1)}{P(\vec{d}|M_2)} \frac{P(M_1)}{P(M_2)}$$

The evidence is larger for a model if more of its parameter space is likely and smaller for a model with large areas in its parameter space having low likelihood values, even if the likelihood function is sharply peaked.

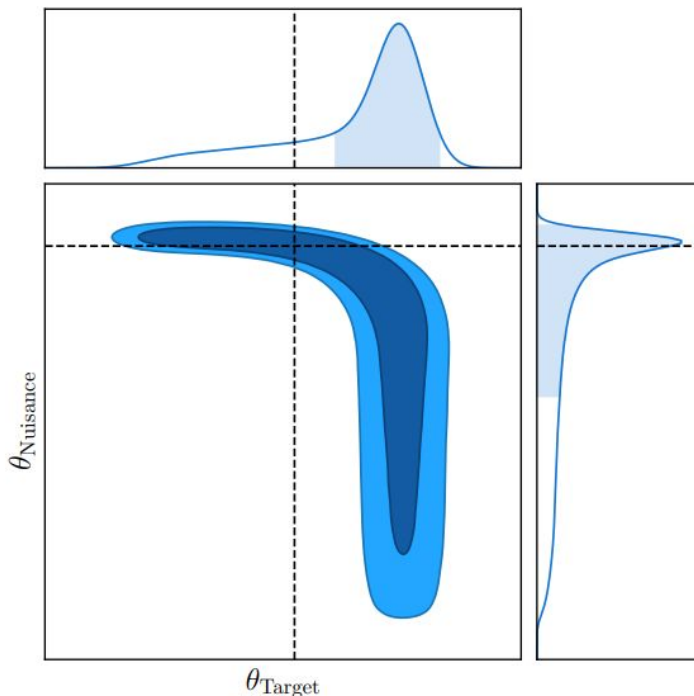
Deviance Information Criterion:

$$\text{DIC}(M) = 2\langle\chi^2\rangle_M - \chi^2_{\text{MaxP}}(M)$$

The model with the lower DIC either fits better the data - lower $\langle\chi^2\rangle$ - or has a lower level of complexity - lower χ^2_{MaxP} . It can be easily computed directly from the parameter posteriors

PRIOR VOLUME EFFECTS

The high-dimensionality of parameter spaces reduces the interpretability of posteriors to their one- and two-dimensional marginal distributions, when more information is available in the full dimensional distributions.



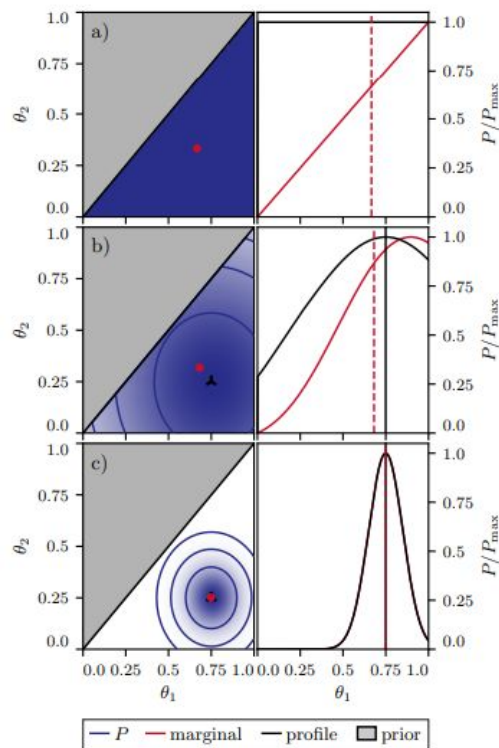
E.g.: Assume a highly-non Gaussian posterior. If we only have access to the 1D marginalized posterior of θ_{Target} , we would hardly estimate the fiducial value with our point estimators (e.g. using the BF point, or the mean/median of the distribution). In other words, even if the posterior in the entire parameter space is centered around the correct parameters, the posterior marginalized over the nuisance parameters can be off of the correct target parameters. This bias error is called the **prior volume or projection effect**.

See e.g.: <https://arxiv.org/pdf/2405.00261>

PRIOR VOLUME EFFECTS

The high-dimensionality of parameter spaces reduces the interpretability of posteriors to their one- and two-dimensional marginal distributions, when more information is available in the full dimensional distributions.

Difference between likelihood profile and marginal distribution for increasingly tighter posterior distributions



Another example is a poorly constrained, prior-limited, parameters, which are projected over significant anisotropic volumes.

→ Posterior profiles do not suffer from projection effects as they are essentially insensitive to the volume of the parameter space. Using a simple metaphor, profiling can be thought of as observing the outline of the posterior landscape, whereas marginalization can be seen as measuring its column density.

<https://arxiv.org/pdf/2409.09101>