

Bayesian Posterior Repartitioning for Nested Sampling

Xi Chen^{*,†}, Farhan Feroz[‡], and Michael Hobson[§]

Abstract. Priors in Bayesian analyses often encode informative domain knowledge that can be useful in making the inference process more efficient. Occasionally, however, priors may be unrepresentative of the parameter values for a given dataset, which can result in inefficient parameter space exploration, or even incorrect inferences, particularly for nested sampling (NS) algorithms. Simply broadening the prior in such cases may be inappropriate or impossible in some applications. Hence our previous solution to this problem, known as posterior repartitioning (PR), redefines the prior and likelihood while keeping their product fixed, so that the posterior inferences and evidence estimates remain unchanged, but the efficiency of the NS process is significantly increased. In its most practical form, PR raises the prior to some power β , which is introduced as an auxiliary variable that must be determined on a case-by-case basis, usually by lowering β from unity according to some pre-defined ‘annealing schedule’ until the resulting inferences converge to a consistent solution. Here we present a very simple yet powerful alternative Bayesian approach, in which β is instead treated as a hyperparameter that is inferred from the data alongside the original parameters of the problem, and then marginalised over to obtain the final inference. We show through numerical examples that this Bayesian PR (BPR) method provides a very robust, self-adapting and computationally efficient ‘hands-off’ solution to the problem of unrepresentative priors in Bayesian inference using NS. Moreover, unlike the original PR method, we show that even for representative priors BPR has a negligible computational overhead relative to standard nesting sampling, which suggests that it should be used as the default in all NS analyses.

Keywords: Bayesian inference, automatic posterior repartitioning, nested sampling, unrepresentative prior.

1 Introduction

In recent years, Monte Carlo sampling techniques have been widely used in Bayesian inference problems both for parameter estimation and model selection. Nested sampling (NS) (Skilling, 2006) is one such approach that can simultaneously produce samples from the posterior distribution and estimate the marginal likelihood (or evidence). The NS algorithm involves drawing samples from a pre-defined prior distribution, and then evaluating their corresponding likelihoods, which are dependent on some measurement (or forward) model and the observed data. In general, the observed data are fixed

*Department of Computer Science, University of Bath, UK. BA2 7PB, xc841@bath.ac.uk

†Department of Physics, University of Cambridge, UK. CB3 0HE, xc253@mrao.cam.ac.uk

‡Department of Physics, University of Cambridge, UK. CB3 0HE, f.feroz@mrao.cam.ac.uk

§Department of Physics, University of Cambridge, UK. CB3 0HE, mph@mrao.cam.ac.uk

and the measurement model is defined by field experts, with little room for flexibility. In contrast, the prior that identifies the regions of interest in the parameter space is typically much more loosely determined and often defined using simple standard distributions, together occasionally with physical constraints on the parameters θ of the problem under consideration.

As we discussed in Chen et al. (2018), the NS algorithm can become very inefficient, or even fail completely in extreme cases, if the likelihood for a given data set is concentrated far out in the wings of the assumed prior distribution. This problem can be particularly damaging in applications where one wishes to perform analyses on many thousands (or even millions) of different datasets, since those (typically few) datasets for which the prior is unrepresentative can absorb a large fraction of the computational resources.

The problem occurs because the NS algorithm begins by drawing a number of ‘live’ samples from the prior and at each subsequent iteration replaces the sample having the lowest likelihood with a sample again drawn from the prior but constrained to have a higher likelihood. Thus, as the iterations progress, the collection of ‘live points’ gradually migrates from the prior to regions of high likelihood. When the likelihood is concentrated very far out in the wings of the prior, this process can become very slow, even in the rare problems where one is able to draw each new sample from the constrained prior using standard methods (sometimes termed perfect nested sampling). In practical problems, the issue is yet more pronounced since algorithms such as MultiNest (Feroz et al., 2009) and PolyChord (Handley et al., 2015) use other methods that may require several likelihood evaluations before a new sample is accepted. Depending on the method used, an unrepresentative prior can thus result in a significant drop in sampling efficiency, thereby increasing still further the required number of likelihood evaluations.

In extreme cases, the migration of live points in the NS process can be very slow, since over many iterations the live points will typically all lie in a region over which the likelihood is very small and flat. Indeed, in some such cases, the log-likelihoods of the set of live points may be indistinguishable to machine precision, so the ‘lowest likelihood’ sample to be discarded will be chosen effectively at random and, in seeking a replacement sample that is drawn from the prior but having a larger likelihood, the algorithm is very unlikely to obtain a sample for which the likelihood value is genuinely larger to machine precision. These problems can result in the live point set migrating exceptionally slowly, or becoming essentially stuck, such that the algorithm (erroneously) terminates before reaching the main body of the likelihood and therefore fails to produce correct posterior samples or evidence estimates.

One may, of course, seek to improve the performance of NS in such cases by increasing the number of live points and/or adjusting the convergence criterion, so that many more NS iterations are performed, but there is no guarantee in any given problem that these measures will be sufficient to prevent premature convergence. Perhaps more useful is to ensure that there is a greater opportunity at each NS iteration of drawing candidate replacement points from regions of the parameter space where the likelihood is larger. This may be achieved in a variety of ways. In MultiNest, for example, one may reduce the `efr` parameter to enlarge the volume of the multi-ellipsoidal bound from which

candidate replacement points are drawn. Alternatively, as in other NS implementations, one may draw candidate replacement points using either MCMC sampling (Feroz and Hobson, 2008) or slice-sampling (Handley et al., 2015) and increase the number of steps taken before a candidate point is chosen. All of these approaches may mitigate the problem to some degree in particular cases, but only at the cost of a simultaneous dramatic drop in sampling efficiency caused precisely by the changes made in obtaining candidate replacement points. Moreover, in more extreme cases, these measures may fail completely.

Aside from making changes to the implementation of the NS algorithm itself, one might consider simple solutions such as broadening the prior range in such cases, but this might not be appropriate or possible in real-world applications, for example when one wishes to assume a single standardised prior across the analysis of a large number of datasets for which the true values of the parameters of interest may vary.

In Chen et al. (2018), we therefore proposed a posterior repartitioning (PR) method that circumvents the above difficulties. The PR method exploits the intrinsic degeneracy between the ‘effective’ likelihood and prior in the formulation of Bayesian inference problems. This is especially relevant for NS since it differs from other sampling methods by making use of the likelihood $\mathcal{L}(\boldsymbol{\theta})$ and prior $\pi(\boldsymbol{\theta})$ *separately* in its exploration of the parameter space, in that samples are drawn from the prior $\pi(\boldsymbol{\theta})$ such that they satisfy $\mathcal{L}(\boldsymbol{\theta}) > L_*$, where L_* is some predefined likelihood constraint. By contrast, Markov chain Monte Carlo (MCMC) sampling methods or genetic algorithm variants are typically blind to this separation,¹ and deal solely in terms of the product $\mathcal{L}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})$, which is proportional to the posterior $\mathcal{P}(\boldsymbol{\theta})$. This difference provides an opportunity in the case of NS to ‘repartition’ the product $\mathcal{L}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})$ by defining a new effective likelihood $\tilde{\mathcal{L}}(\boldsymbol{\theta})$ and prior $\tilde{\pi}(\boldsymbol{\theta})$ (which is typically ‘broader’ than the original prior), subject to the condition $\tilde{\mathcal{L}}(\boldsymbol{\theta})\tilde{\pi}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})$, so that the (unnormalised) posterior remains unchanged. Thus, in principle, the inferences obtained are unaffected by the use of the PR method, but, as we demonstrated in Chen et al. (2018), the approach can yield significant improvements in sampling efficiency and also helps to avoid the convergence problems that can occur in extreme examples of unrepresentative priors.

In general, the effective prior $\tilde{\pi}(\boldsymbol{\theta})$ can be any distribution that can be sampled straightforwardly, which in principle can be arbitrary (Alsing and Handley, 2021), but should at least be non-zero everywhere that the original prior is non-zero. In its most practical form, however, sometimes termed power posterior repartitioning (PPR),² the effective prior $\tilde{\pi}(\boldsymbol{\theta})$ is proportional simply to the original prior $\pi(\boldsymbol{\theta})$ raised to some power β , which is introduced as an auxiliary variable. Although we demonstrated the effectiveness of this approach in Chen et al. (2018), one drawback of our original method is that the auxiliary variable must be determined on a case-by-case basis, which is typically achieved by lowering β from unity, outside the execution of the NS algorithm, according

¹One exception is the propagation of multiple MCMC chains, for which it is often advantageous to draw the starting point of each chain independently from the prior distribution.

²This naming convention, which we shall avoid here, should not be confused with the power posterior method (also known as thermodynamic integration) for estimating the evidence using MCMC sampling; the latter instead involves transitioning from the prior to the posterior by powering just the likelihood by an inverse temperature.

to some pre-defined ‘annealing schedule’, until the resulting inferences from successive NS runs converge to a statistically consistent solution for values below some (positive) threshold $\beta \lesssim \beta_*$. This approach lacks elegance and the repeated NS runs required can be computationally demanding. Moreover, the final inference is unavoidably conditioned on the value β_* .

In this paper, we therefore present an alternative, Bayesian approach, in which β is instead treated as a hyperparameter that is inferred from the data alongside the original parameters $\boldsymbol{\theta}$ of the problem, within a single execution of the NS algorithm. Although this approach is conceptually very simple, indeed almost trivial, the resulting Bayesian PR (BPR) method is considerably more powerful than our original PR technique in several respects. In particular, one obtains samples from the joint posterior on $(\boldsymbol{\theta}, \beta)$, which may then be used straightforwardly in a number of ways. First, one may properly marginalise over β in a Bayesian manner to obtain a final inference on $\boldsymbol{\theta}$, rather than conditioning on a single value of β as in the original PR approach. Moreover, this allows BPR to accommodate likelihood functions consisting of a number of spatially separated modes that are located asymmetrically with respect to the prior, and are therefore characterised by different ranges of β values; this is not possible using the original PR method. Second, one may instead marginalise over $\boldsymbol{\theta}$ to determine the 1-dimensional marginal distribution of β (an ‘effective’ posterior that will be introduced in later sections), which we will show is very useful in diagnosing both the existence and severity of an unrepresentative prior for a given dataset, and hence for identifying ‘outlier’ datasets and in refining the inference problem for future analyses. Finally, since the sampling of the joint posterior on $(\boldsymbol{\theta}, \beta)$ is performed within a single NS run, BPR is much less computationally demanding than the original PR method, which requires a separate NS run for each value of β in its annealing schedule. Indeed, since the overhead of introducing just a single additional hyperparameter β to be sampled is negligible for most practical problems, BPR is typically no more computationally demanding than a standard NS analysis (i.e., equivalent to setting $\beta = 1$), but automatically safeguards against potentially inefficient parameter space exploration, or even incorrect inferences, which may occur in the presence of unrepresentative priors. This suggests that BPR should, in fact, be used as the default approach in all NS analyses.

It is worth noting, however, that one may encounter even more extreme problems than those discussed above, where the likelihood for some dataset(s) is concentrated outside an assumed prior having compact support. From a probabilistic point of view, the prior in these problems has zero probability in the region of the likelihood distribution. This case, which one might describe as an unsuitable prior, is not addressed by the BPR method, and is not considered here.

This paper is organised as follows. Section 2 briefly introduces the Bayesian inference and the NS algorithm. In Section 3, we present a simple analytical illustration of the effect of an unrepresentative prior on Monte Carlo sampling in general and NS in particular. Section 4 outlines the original PR method and then describes the proposed BPR scheme. Section 5 demonstrates performance of the BPR method in some numerical examples, and we present our conclusions in Section 6.

2 Bayesian inference using nested sampling

Bayesian inference (see e.g. MacKay 2003) provides a consistent framework for estimating unknown parameters θ of some model by updating any prior knowledge of θ using the observed data \mathcal{D} and an assumed measurement process. The complete inference is embodied in the posterior distribution of θ , which can be expressed using Bayes' theorem as:

$$\Pr(\theta|\mathcal{D}, \mathcal{M}) = \frac{\Pr(\mathcal{D}|\theta, \mathcal{M}) \Pr(\theta|\mathcal{M})}{\Pr(\mathcal{D}|\mathcal{M})}, \quad (1)$$

where \mathcal{M} represents model (or hypothesis) assumption(s), and adopting a simplifying notation, $\Pr(\theta|\mathcal{D}, \mathcal{M}) \equiv \mathcal{P}(\theta)$ is the *posterior* probability density, $\Pr(\mathcal{D}|\theta, \mathcal{M}) \equiv \mathcal{L}(\theta)$ is the *likelihood*, $\Pr(\theta|\mathcal{M}) \equiv \pi(\theta)$ is the *prior* probability density on θ and $\Pr(\mathcal{D}|\mathcal{M}) \equiv \mathcal{Z}$ is called the *evidence* (or marginal likelihood). We then have the simplified expression:

$$\mathcal{P}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{\mathcal{Z}}, \quad (2)$$

in which the evidence is given by

$$\mathcal{Z} = \int_{\mathcal{R}} \mathcal{L}(\theta)\pi(\theta)d\theta, \quad (3)$$

where \mathcal{R} represents the prior space of θ . The evidence \mathcal{Z} is often used for model selection. It is the average of the likelihood over the prior, considering every possible choice of θ , and thus is not a function of the parameters θ . The constant \mathcal{Z} is usually ignored in parameter estimation, since the posterior $\mathcal{P}(\theta)$ is proportional to the product of likelihood $\mathcal{L}(\theta)$ and prior $\pi(\theta)$.

Nested sampling (Skilling, 2006) explores the posterior distribution in a sequential manner using a fixed number of N_{live} 'live samples' of the parameters θ at each iteration of the process.³ Pseudo-code of standard NS algorithm is described in Supplementary Materials (Chen et al., 2022). Among the various implementations of the NS algorithm (Chopin and Robert, 2007; Feroz et al., 2009; Brewer et al., 2011; Handley et al., 2015; Baldock et al., 2017; Higson et al., 2018; Buchner, 2019; Speagle, 2020; Buchner, 2021), two widely used packages are MultiNest (Feroz et al., 2009, 2013) and PolyChord (Handley et al., 2015). MultiNest draws the new sample at each iteration using rejection sampling from within a multi-ellipsoid bound approximation to the iso-likelihood surface defined by the discarded point; the bound is constructed from the live points present at that iteration. PolyChord draws the new sample at each iteration using a number of successive slice-sampling steps taken in random directions, which is particularly well suited to higher-dimensional problems. Please see Feroz et al. (2009) and Handley et al. (2015) for more details.

³The method has recently been extended to so-called dynamic nested sampling (Higson et al., 2018), which allows the number of live samples to vary as the NS iterations proceed, but we will not use this variant here.

3 Analytical illustration of Monte Carlo sampling with an unrepresentative prior

We begin by considering a simple intuitive example that illustrates the effect of an unrepresentative prior on Monte Carlo sampling algorithms in general and on NS in particular. We consider the case where both the likelihood and prior distributions are one-dimensional Gaussians defined by $\mathcal{L}(\theta) = \mathcal{N}(\theta; 0, 1)$ and $\pi(\theta) = \mathcal{N}(\theta; \mu_\pi, 1)$, respectively, where μ_π is the mean of prior and θ is rewritten as θ for the one dimensional example. In this case, the fractional prior volume X (which ranges between $[0, 1]$ for NS) occupied by the region for which the likelihood exceeds the level $\mathcal{L} = \lambda$ is given by

$$X(\lambda) = \int_{-|\theta(\lambda)|}^{|\theta(\lambda)|} \pi(\theta) d\theta = \frac{1}{2} \left[\operatorname{erf} \left(\frac{|\theta(\lambda)| - \mu_\pi}{\sqrt{2}} \right) - \operatorname{erf} \left(-\frac{|\theta(\lambda)| + \mu_\pi}{\sqrt{2}} \right) \right], \quad (4)$$

where erf is the error function and $\theta(\lambda) = \pm \sqrt{-\log(2\pi) - 2\log(\lambda)}$ is the value of random variable θ where the likelihood value is λ . Also, the fractional evidence contained in the complement region where the likelihood is below the level $\mathcal{L} = \lambda$ is:

$$\begin{aligned} \mathcal{Z}(X = X(\lambda)) &= \left(\int_{-\infty}^{\infty} - \int_{-|\theta(\lambda)|}^{|\theta(\lambda)|} \right) \mathcal{L}(\theta) \pi(\theta) d\theta \\ &= \frac{1}{2\sqrt{\pi}} \exp \left(-\frac{\mu_\pi^2}{4} \right) \left[1 - \frac{1}{2} \left\{ \operatorname{erf} \left(\frac{\mu_\pi}{2} + |\theta(\lambda)| \right) - \operatorname{erf} \left(\frac{\mu_\pi}{2} - |\theta(\lambda)| \right) \right\} \right]. \end{aligned} \quad (5)$$

We consider three cases for which the mean μ_π of the Gaussian prior equal to 1, 5 and 15, respectively. In the first case, the prior and likelihood are consistent, while the other two cases correspond to increasingly unrepresentative priors. Figure 1 shows (a) the behaviour of the likelihood \mathcal{L} and (b) the evidence \mathcal{Z} as a function of the prior volume X , all on a logarithmic scale.

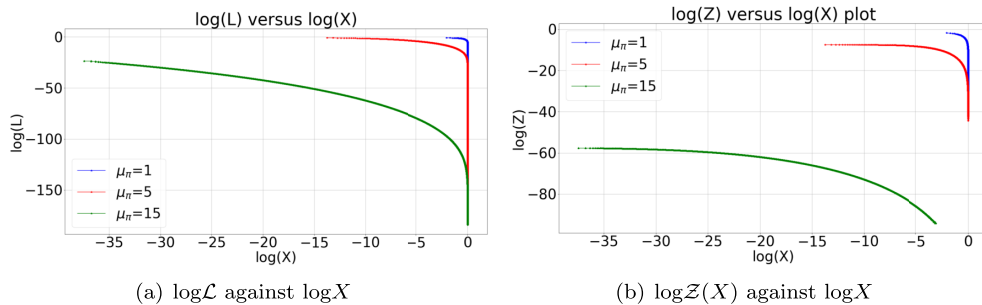


Figure 1: Relationship between likelihood $\mathcal{L}(X)$ and evidence $\mathcal{Z}(X)$ (up to prior volume X) against the prior volume X for the toy examples discussed in Section 3. The three curves in each figure represent the priors with mean μ_π equal to 1, 5 and 15, respectively.

This figure shows that both $\mathcal{L}(X) \approx 0$ and $\mathcal{Z}(X) \approx 0$ (i.e., large negative values on a log-scale) in the case of an extreme unrepresentative with $\mu = 15$ (green curve) for all values of X except those very close to zero (or $\log X \approx -\infty$). Also both $\log \mathcal{L}(X)$ and $\log \mathcal{Z}(X)$ increase more gradually with decreasing $\log X$ as μ_π increases (and the prior becomes more unrepresentative). Therefore, for higher μ_π it becomes increasingly difficult for any Monte Carlo algorithm in general and NS in particular to proceed as the likelihood is almost zero over most of the prior space and increases very gradually as the NS algorithm proceeds to regions contained within iso-likelihood contours at progressively higher likelihood levels, resulting in premature convergence. This also explains the results shown in Figure 13(a) in the Supplementary Material, where standard NS experiences a sudden failure as the prior becomes increasingly unrepresentative.

4 Bayesian posterior repartitioning

In previous research, we first considered the possibility of improving the robustness and efficiency of NS by exploiting the intrinsic degeneracy between the ‘effective’ likelihood and prior in the formulation of Bayesian inference problems in Feroz et al. (2009), and then further developed the idea as the original PR method in Chen et al. (2018).

As mentioned in the Introduction, the central idea is to ‘repartition’ the product of the likelihood and prior, such that

$$\mathcal{L}(\boldsymbol{\theta})\pi(\boldsymbol{\theta}) = \tilde{\mathcal{L}}(\boldsymbol{\theta})\tilde{\pi}(\boldsymbol{\theta}), \quad (6)$$

where $\tilde{\mathcal{L}}(\boldsymbol{\theta})$ and $\tilde{\pi}(\boldsymbol{\theta})$ are the newly-defined effective (or modified) likelihood and prior, respectively. As a result, the (unnormalised) posterior remains unchanged and hence, in principle, any parameter inferences are unaffected. Moreover, if $\tilde{\pi}(\boldsymbol{\theta})$ is normalised, then the evidence also remains unchanged. In general, $\tilde{\pi}(\boldsymbol{\theta})$ may be any distribution that can be straightforwardly sampled, and in principle can be arbitrary with sufficient effort (Alsing and Handley, 2021). For example, there is no requirement for the modified prior to be centred at the same parameter value as the original prior. One could, therefore, choose a modified prior that broadens and/or shifts the original one, or a modified prior that has a completely different form from the original. In this generalised setting, however, the modified prior should at least have the same support as the original prior.

In practice, rather than introducing a completely new prior distribution into the problem, a convenient choice (which we previously termed power PR) is simply to take $\tilde{\pi}(\boldsymbol{\theta})$ to be the original prior $\pi(\boldsymbol{\theta})$ raised to some (real) power β , and then renormalised to unit volume, such that

$$\tilde{\pi}(\boldsymbol{\theta}) = \frac{\pi(\boldsymbol{\theta})^\beta}{\mathcal{Z}_\pi(\beta)}, \quad (7)$$

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})^{(1-\beta)}\mathcal{Z}_\pi(\beta), \quad (8)$$

where $\beta \in [0, 1]$ and $\mathcal{Z}_\pi(\beta) \equiv \int \pi(\boldsymbol{\theta})^\beta d\boldsymbol{\theta}$ is the normalisation constant of the modified prior. By altering the value of β , the modified prior varies across a range of distributions between the original prior ($\beta = 1$) and the uniform distribution ($\beta = 0$). This is

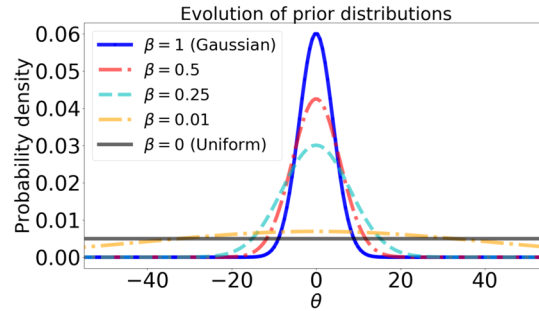


Figure 2: One-dimensional prior evolution for $\beta \in [0, 1]$. The original prior is a Gaussian distribution with $\sigma_\pi = 4$ (truncated in the range $[-50, 50]$) when $\beta = 1$ (solid blue curve), and is an uniform distribution when $\beta = 0$ (solid black curve). The remaining three curves correspond to $\beta = 0.5$ (dot-dashed red curve), 0.25 (dashed light blue curve), 0.01 (dot-dashed yellow curve), respectively.

illustrated in Figure 2 for five specific β values in a one-dimensional case, where the original prior is a Gaussian with zero mean and standard deviation $\sigma_\pi = 4$ and the normalisation depends on the assumed support $[-50, 50]$ of the unknown parameter θ . The $\beta = 0$ limit clearly yields a uniform modified prior $\tilde{\pi}(\theta) \sim \mathcal{U}(\mathcal{R})$ over the allowed region \mathcal{R} of the parameter space, which is an important special case in that it maximises the dispersion of the initial live point set across the prior space, but consequently can result in very inefficient sampling. It is worth noting that in the special case in which the original prior is uniform $\pi(\theta) \sim \mathcal{U}(\mathcal{R})$, the power PR method is insensitive to the value of β and defaults to standard NS. In principle, one may extend the upper limit on β to exceed unity; this produces a new prior distribution that is sharper than the original, which may prove useful if the latter is over-dispersed. The power PR method does, however, have the limitation that one must be able to sample from the resulting modified prior and evaluate the normalising constant $\mathcal{Z}_\pi(\beta)$. If necessary, the latter may be estimated numerically, most efficiently in a separate NS run, for each value of β used, but this may increase the computational costs significantly. Nonetheless, this issue may potentially be mitigated by pre-calculating $\mathcal{Z}_\pi(\beta)$ for a set of β values and storing them in a look-up table, from which interpolation can be performed to approximate $\mathcal{Z}_\pi(\beta)$ for other β values.

A clear drawback of the original PR method, however, is that an appropriate value of the auxiliary variable β must be determined on a case-by-case basis, and this can depend sensitively on the nature of the original prior and likelihood, as well as on the dimensionality of the problem under consideration. In Chen et al. (2018), we therefore suggested an approach in which β is gradually lowered from unity (outside of the NS algorithm) according to some ‘annealing schedule’, until the resulting inferences from successive NS runs converge to a statistically consistent solution, which typically occurs for β values below some (positive) threshold, $\beta \lesssim \beta_*$. Although we demonstrated in Chen et al. (2018) that this method is robust and effective, it has the disadvantages that there is a significant computational overhead associated with the multiple NS runs

required and that the final inference of the parameters of interest $\boldsymbol{\theta}$ is conditioned on the adopted value $\beta = \beta_*$.

We therefore propose an alternative approach here, which we consider to be considerably more elegant, whereby β is instead treated as a hyperparameter that is inferred in a fully Bayesian manner alongside the original parameters $\boldsymbol{\theta}$, within a single run of the NS algorithm. Although this approach is admittedly conceptually very simple, we view this as a virtue, in particular because the resulting Bayesian PR (BPR) method is considerably more powerful and flexible than the original PR technique in a number of important ways. The BPR method is based on defining the joint posterior

$$\tilde{\mathcal{P}}(\boldsymbol{\theta}, \beta) \propto \tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta) \tilde{\pi}(\boldsymbol{\theta}, \beta) = \tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta) \tilde{\pi}(\boldsymbol{\theta}|\beta) \pi(\beta), \quad (9)$$

where $\pi(\beta)$ denotes the assumed prior on the hyperparameter β , and $\tilde{\pi}(\boldsymbol{\theta}|\beta)$ and $\tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta)$ have precisely the forms (7) and (8), respectively. Since β lies naturally in the range $[0, 1]$, we define $\pi(\beta)$ to be the uniform prior over this interval, although other choices may be accommodated if there is a strong motivation to adopt an alternative form in a particular problem. For example, if one wishes to extend the upper limit on β to exceed unity to accommodate an over-dispersed original prior $\pi(\boldsymbol{\theta})$ on the parameters, a natural prior on β is $\pi(\beta) = e^{-\beta}$ in the range $[0, \infty]$, which is the maximum-entropy distribution for a non-negative quantity for which one knows only that its expectation value is unity. In any case, from the relations (6) and (2), one sees that (by construction)

$$\tilde{\mathcal{P}}(\boldsymbol{\theta}, \beta) \propto \mathcal{P}(\boldsymbol{\theta}) \pi(\beta). \quad (10)$$

Moreover, the corresponding evidence is given by

$$\tilde{\mathcal{Z}} = \iint \tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta) \tilde{\pi}(\boldsymbol{\theta}, \beta) d\boldsymbol{\theta} d\beta \quad (11)$$

$$= \int \mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \int \pi(\beta) d\beta = \mathcal{Z}, \quad (12)$$

so that the proportionality in (10) can be replaced by an equality.

In principle, a single NS run thus provides samples from the full joint posterior distribution (10), which can be used straightforwardly to obtain inferences on the original parameters $\boldsymbol{\theta}$ in a properly Bayesian manner by marginalising over β , since

$$\tilde{\mathcal{P}}(\boldsymbol{\theta}) = \int \tilde{\mathcal{P}}(\boldsymbol{\theta}, \beta) d\beta = \int \mathcal{P}(\boldsymbol{\theta}) \pi(\beta) d\beta = \mathcal{P}(\boldsymbol{\theta}). \quad (13)$$

Hence, the overall computational burden is much reduced compared to the original PR method, since one does not require the latter's multiple NS runs, namely one for each value of β in its annealing schedule, or more if runs with the same value of β (particularly for large β values) are liable to exhibit significant variation in their results because of the failure of the NS algorithm, thereby complicating the process of determining convergence of results as the annealing schedule proceeds. Moreover, the final inference on the parameters $\boldsymbol{\theta}$ in BPR is not conditioned on a particular value of β . As we will show in examples in Section 5, marginalising over β also allows BPR to accommodate

likelihood functions consisting of multiple spatially separated modes that are located asymmetrically with respect to the prior, and which are therefore characterised by different ranges of β values; this is not possible using the original PR method.

A further advantage of BPR is that one may instead choose to marginalise over the original parameters $\boldsymbol{\theta}$ to obtain the 1-D ‘effective’ posterior on β . This should in principle simply recover the prior on β , since

$$\tilde{\mathcal{P}}(\beta) = \int \tilde{\mathcal{P}}(\boldsymbol{\theta}, \beta) d\boldsymbol{\theta} = \int \mathcal{P}(\boldsymbol{\theta})\pi(\beta) d\boldsymbol{\theta} = \pi(\beta). \quad (14)$$

Thus, if the NS procedure has correctly sampled from the joint posterior distribution $\tilde{\mathcal{P}}(\boldsymbol{\theta}, \beta)$, one has neither gained nor lost anything by introducing the hyperparameter β , apart perhaps from a slight increase in the computational burden as a result of increasing by one the dimensionality of the space to be sampled.

In practice, however, the situation is more subtle. For illustration, let us consider the case where the original prior $\pi(\boldsymbol{\theta}) = \tilde{\pi}(\boldsymbol{\theta}, 1)$ is extremely unrepresentative of the dataset under analysis, so that the original likelihood $\mathcal{L}(\boldsymbol{\theta}) = \tilde{\mathcal{L}}(\boldsymbol{\theta}, 1)$ is concentrated very far into the wings of $\pi(\boldsymbol{\theta})$. In the limit $N_{\text{live}} \rightarrow \infty$, the NS algorithm would nonetheless converge correctly, yielding samples from the posterior (10) and an estimate of the evidence (12). For finite N_{live} , however, live points drawn from $\tilde{\pi}(\boldsymbol{\theta}, \beta)$ at any NS iteration will typically have very low likelihood $\tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta)$ for values of β above some limiting threshold $\beta \gtrsim \beta_+$ (which will depend on N_{live}), since the chance of these points lying within the main body of the likelihood $\tilde{\mathcal{L}}(\boldsymbol{\theta}, \beta)$ is vanishingly small. Depending on the precise nature of the original prior $\pi(\boldsymbol{\theta})$ and likelihood $\mathcal{L}(\boldsymbol{\theta})$, a similar phenomenon may also occur for values of β below some other limiting threshold $\beta \lesssim \beta_-$ (which will also depend on N_{live}), since in this case the modified prior may have a support that far exceeds that of the likelihood.

Thus, in the presence of unrepresentative priors where the NS algorithm may become very inefficient or even fail, one may obtain samples that are drawn not from (10), but instead from some ‘effective’ posterior

$$\tilde{\mathcal{P}}_{\text{eff}}(\boldsymbol{\theta}, \beta) \propto \mathcal{P}(\boldsymbol{\theta})\tilde{\mathcal{P}}(\beta), \quad (15)$$

where the (unnormalised) marginal posterior $\tilde{\mathcal{P}}(\beta)$ is non-zero only in some range $\sim [\beta_-, \beta_+]$. The form of this marginal posterior, in particular its extent in β , will be determined by how the NS algorithm fails for particular (extreme) values of β , which is in turn dependent on the particular NS implementation used and the associated control parameters (including generic NS parameters such as N_{live}), as well as the nature of the original prior $\pi(\boldsymbol{\theta})$ and likelihood $\mathcal{L}(\boldsymbol{\theta})$. One would, however, expect there to be a range of β values for which the NS algorithm does not fail, and so $\tilde{\mathcal{P}}(\beta)$ should coincide with $\pi(\beta)$ in this range. Beyond these very general considerations, there is a paucity of further theoretical arguments from which to predict the form of $\tilde{\mathcal{P}}(\beta)$.

Nonetheless, by marginalising (15) over β one still obtains the posterior $\mathcal{P}(\boldsymbol{\theta})$ on the original parameters. Conversely, marginalising over $\boldsymbol{\theta}$ one obtains $\tilde{\mathcal{P}}(\beta)$ (in practice as a histogram constructed from the equally weighted posterior samples) and may hence

observe its form and determine the values β_- and β_+ . In particular, the value of β_+ is useful in diagnosing both the existence and severity of an unrepresentative prior for a given dataset, and thereby identifying the dataset as an ‘outlier’. Here we will take β_- and β_+ simply as the smallest and largest β values, respectively, in the set of equally weighted posterior samples. Alternatively, one could apply some percentile thresholds (e.g. 1% and 99%) to the β marginal, but in practice this leads to very similar values of β_- and β_+ . In any case, observing the form of $\tilde{\mathcal{P}}(\beta)$ obtained in the BPR method provides a far more robust and computationally efficient approach to identifying the presence and severity of unrepresentative priors than determining and interpreting the value β_* in the original PR method below which the inferences converge.

Finally, in the presence of unrepresentative priors, the NS process will not yield an estimate of the evidence (12), but rather the ‘effective’ evidence

$$\tilde{\mathcal{Z}}_{\text{eff}} \approx \tilde{\mathcal{Z}} \int \tilde{\mathcal{P}}(\beta) d\beta = \mathcal{Z} \int \tilde{\mathcal{P}}(\beta) d\beta. \quad (16)$$

One may, however, estimate the required evidence \mathcal{Z} , provided one can evaluate the factor $\int \tilde{\mathcal{P}}(\beta) d\beta$. Fortunately, from the discussion above, this may be achieved by first scaling the distribution $\tilde{\mathcal{P}}(\beta)$, represented as a histogram of equally-weighted posterior samples, such that the bin containing the largest number of samples has the volume $\int_{\beta_l}^{\beta_u} \pi(\beta) d\beta$, where β_l and β_u are the lower and upper limits of the bin, respectively. The factor $\int \tilde{\mathcal{P}}(\beta) d\beta$ may then be calculated by summing up the volumes of all the bins in the resulting scaled histogram. This summation is equivalent to the simplest form of numerical quadrature, namely the rectangle rule, which adopts a polynomial of degree zero (a constant function) to pass through points $((\beta_l^b + \beta_u^b)/2, h(\beta_l^b + \beta_u^b)/2)$, where β_l^b and β_u^b are the lower and upper limits of the b th bin, and $h(\cdot)$ indicates the height of the bin. In principle, one may estimate the error in the quadrature and propagate this through to $\tilde{\mathcal{Z}}_{\text{eff}}$ using (16), but we do not perform this here.

5 Numerical examples

We now illustrate the performance of the BPR method in some numerical examples. Since we extensively studied the behaviour of our original PR method in a wide range of problems in Chen et al. (2018), including comparisons with other sampling algorithms such as Markov Chain Monte Carlo (MCMC) and importance sampling, we focus here primarily on the performance of BPR in the canonical case where the likelihood and prior on the original variables θ are both Gaussian (although very mismatched in some examples), and hence so too is the posterior. In particular, we begin by considering the univariate case, before moving on to a bivariate example for which we consider both circularly-symmetric and asymmetric priors, where the latter may have zero, positive or negative correlation coefficient, respectively. To demonstrate the method further, however, we also consider several 2-dimensional non-Gaussian likelihoods: multi-modal likelihoods consisting of an equal mixture of four identical but spatially separated Gaussians, and a unimodal Laplacian likelihood (the latter example is presented in the Supplementary Material). Finally, we investigate higher-dimensional unimodal Gaussian

likelihood examples, up to 10 dimensions. In these further examples, we consider only circularly-symmetric Gaussian priors for the sake of brevity.

In all our numerical examples, we use the NS package MultiNest (Feroz et al., 2009), with algorithm parameter settings similar to those used in Chen et al. (2018) for the study of our original PR method. Specifically, we set the number of live points $N_{\text{live}} = 100$, the desired sampling efficiency parameter $\text{efr} = 0.8$, which equals the ratio of the remaining prior volume to the volume of the multi-ellipsoid bound enclosing the active point set at each NS iteration, and the convergence tolerance parameter $\text{tol} = 0.5$, which corresponds to acceptable uncertainty on the estimated log-evidence (please refer to Feroz et al. (2009) for more details).

5.1 Univariate Gaussian likelihood

We begin by considering a simple one-dimensional estimation problem, for which the data consist of N independent measurements $M = \{m_1, \dots, m_n, \dots, m_N\}$ of an unknown parameter θ , such that

$$m_n = \theta + \xi, \quad (17)$$

where ξ denotes Gaussian noise $\xi \sim \mathcal{N}(\mu_\xi, \sigma_\xi^2)$. The likelihood is thus given by

$$\mathcal{L}(\theta) = \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma_\xi^2}} \exp \left[-\frac{(m_n - (\theta + \mu_\xi))^2}{2\sigma_\xi^2} \right] \right\}, \quad (18)$$

which has a Gaussian form. In particular, we set $\mu_\xi = 0$, $\sigma_\xi = 1$ and the number of measurements as $N = 20$, so that the likelihood is a Gaussian centred around the true value of θ with a standard deviation of $\sim 1/\sqrt{20} \approx 0.22$. The prior distribution $\pi(\theta)$ is also assumed to be Gaussian $\theta \sim \mathcal{N}(\mu_\pi, \sigma_\pi^2)$, where we set $\mu_\pi = 0$ and $\sigma_\pi = 4$. Therefore one expects *a priori* that the true value of θ will lie in the range $[-12, 12]$ with greater than 99.5% probability. Finally, we assume the prior distribution of the auxiliary factor β to be the unit uniform distribution $\beta \sim \mathcal{U}[0, 1]$.

We consider datasets corresponding to a series of true values θ_* ranging between 5 and 50. Figure 3(a) shows the limiting example for $\theta_* = 50$, where the likelihood lies very far into the wings of the prior and may thus be considered unrepresentative for this data set. Figure 3(b) shows the posterior samples obtained using standard MultiNest without PR, together with the corresponding true Gaussian posterior. The behaviour in the limiting case is characteristic of a catastrophic failure mode of the NS algorithm in practice for extreme unrepresentative priors. This was first discussed in Chen et al. (2018) and here we take a step further to illustrate and discuss this behaviour in Figure 13 and Section 7.3 of the Supplementary Materials. Moreover, a discussion on failure probability of standard NS is also presented using the same univariate example as follows.

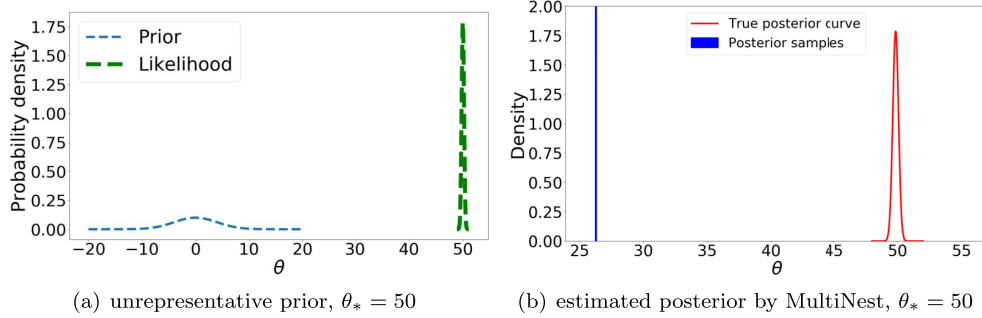


Figure 3: Univariate Gaussian likelihood examples illustrating the unrepresentative case and the posterior samples obtained by standard MultiNest (blue histogram line). The light blue and green dashed curves denote the prior and the likelihood, respectively. The red curve represents the true Gaussian posterior.

Failure probability of NS with univariate Gaussian likelihood

We calculate the failure probability of standard NS for the case where prior is a univariate Gaussian with mean μ_π and standard deviation σ_π . The likelihood is also a univariate Gaussian with mean μ_ξ and standard deviation σ_ξ with $\mu_\xi > \mu_\pi$. We consider failure to be the case where there are no live points within $c\sigma_\xi$ of the likelihood mean μ_ξ after I iterations. One could use $c = 3$ here. The failure will very likely happen before I iterations due to all live points being the same within machine precision, so we can take the failure probability estimate p_{fail} obtained here as its lower bound.

Before the first iteration, for the failure to happen, all N_{live} points have to be at values less than $(\mu_\xi - c\sigma_\xi)$ which has probability:

$$p_{\text{fail},0} = \Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi)^{N_{\text{live}}}, \tag{19}$$

where $\Phi(x; \mu, \sigma)$ is the cumulative distribution function of Normal distribution and $p_{\text{fail},0}$ denotes the failure probability estimate at the $i = 0$ iteration.

At each subsequent iteration i , if the remaining prior volume is v_i and the point with minimum likelihood value is at position y_i , then we have:

$$1 - \Phi(y_i; \mu_\pi, \sigma_\pi) = v_i. \tag{20}$$

We can then calculate the probability of failure at iteration i , $p_{\text{fail},i}$ as:

$$\begin{aligned} p_{\text{fail},i} &= \int_0^1 \frac{P(v_i) (\Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi) - \Phi(y_i; \mu_\pi, \sigma_\pi))}{v_i} dv_i, \\ &= \int_0^1 \frac{P(v_i) (\Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi) - 1 + v_i)}{v_i} dv_i, \end{aligned} \tag{21}$$

where the integral is over the distribution $P(v_i)$ of remaining prior volume v_i at iteration i which itself equals to $v_i = \prod_{j=1}^i t_j$, where $P(t) = N_{\text{live}} t^{N_{\text{live}}-1}$ (Feroz et al., 2009) i.e.

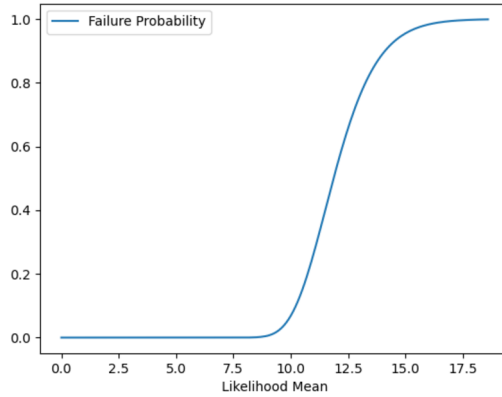


Figure 4: Failure probability of standard NS for the univariate Gaussian example as a function of mean of the likelihood distribution μ_ξ (which is equal to θ_* in the rest of this section).

the distribution of the maximum of N_{live} uniformly distributed random variables and therefore v_i is the product of beta distributed random variables. It isn't possible to get an analytical solution for the integral given in (21). However, we can get an approximation by assuming $v_i \approx \mathbb{E}[v_i]$, where $\mathbb{E}[v_i] = \exp(-i/N_{\text{live}})$ (Feroz et al., 2009). Using this approximation in (21) we get:

$$p_{\text{fail},i} \approx \frac{\Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi) - 1 + \mathbb{E}[v_i]}{\mathbb{E}[v_i]}. \quad (22)$$

We can now approximate the total failure probability p_{fail} as:

$$\begin{aligned} p_{\text{fail}} &\approx p_{\text{fail},0} \prod_{i=1}^I p_{\text{fail},i} \\ &= \Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi)^{N_{\text{live}}} \prod_{i=1}^I \frac{\Phi(\mu_\xi - c\sigma_\xi; \mu_\pi, \sigma_\pi) - 1 + \mathbb{E}[v_i]}{\mathbb{E}[v_i]}. \end{aligned} \quad (23)$$

We can use p_{fail} given above to explain the failure of standard NS in the univariate example discussed in this Section, i.e., $c = 3$, $I = 100$, $N_{\text{live}} = 100$, $\mu_\pi = 0$, $\sigma_\pi = 4$, $\sigma_\xi = 0.22$. We can now calculate the failure probability p_{fail} for different values of μ_ξ (which is equal to θ_* in Figure 5). As shown in Figure 4, the failure probability increases rapidly for $\mu_\xi > 10$ and therefore, we would expect $\beta < 1$ for $\mu_\xi > 10$ with BPR which is indeed what we see in Figure 5.

The above illustrated issue of NS may be straightforwardly addressed by applying the BPR method. Figure 5 shows the resulting MultiNest (with BPR) joint posterior on (θ, β) and its marginals for a selection of true values θ_* in the range $[5, 50]$ (the $\theta_* = 50$ case is plotted in brown). One sees that the joint posterior is precisely of the

form expected in (15), in that it is the product of two independent distributions on θ and β , respectively. One also observes the evolution of the marginal distribution on β . For each value of θ_* , this is consistent with a top-hat distribution in the range $[0, \beta_+]$, where β_+ gradually decreases as θ_* increases. For $\theta_* \lesssim 10$, one sees that $\beta_+ \approx 1$, so one recovers the original uniform prior distribution $\beta \sim \mathcal{U}[0, 1]$, indicating that the original prior $\pi(\theta)$ is representative for these data sets. For $\theta_* \gtrsim 10$, however, the value of β_+ gradually decreases from unity as θ_* increases, which indicates that the original prior is increasingly unrepresentative for these data sets.

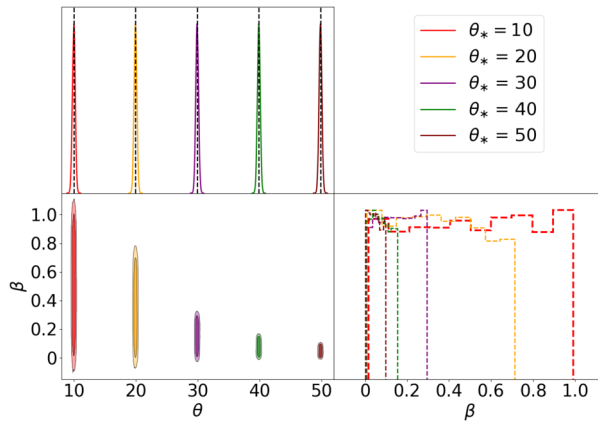


Figure 5: A ‘corner plot’ showing the joint posterior distribution on (θ, β) obtained by MultiNest using BPR, together with its marginals, for the univariate Gaussian likelihood example with a range of true values θ_* . The vertical dashed black lines in the top left panel indicate the mean of the true Gaussian posterior in each case.

Turning to the marginals on θ in Figure 5, one sees that they are correctly centred on the mean of the corresponding true Gaussian posterior for every value of θ_* . Moreover, the widths of the marginals on θ are equal for each value of θ_* and consistent with the width of the true Gaussian posterior, hence showing that BPR yields the correct inferences on θ , independent of the value of θ_* , in a fully automated manner, without any need for tuning. Quantitative results on both the parameter estimation accuracy on θ and the evidence calculation are given in the Supplementary Material Section 7.3.

5.2 Bivariate Gaussian likelihood

We now move to a bivariate example, for which we consider both circularly-symmetric and asymmetric priors, where the latter may have zero, positive or negative correlation coefficient, respectively. In particular, consider a vectorised version of (17) from the univariate Gaussian likelihood example, with some $K = 2$ dimensional parameter vector $\theta = (\theta_1, \theta_2)^\top$, such that

$$\mathbf{m}_n = \theta + \xi, \tag{24}$$

where $\mathbf{m} = (m_1, m_2)^\top$, and $\boldsymbol{\xi} = (\xi_1, \xi_2)^\top$ denotes the two dimensional Gaussian noise $\boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}_\xi, \boldsymbol{\Sigma}_\xi)$ with mean $\boldsymbol{\mu}_\xi$ and covariance $\boldsymbol{\Sigma}_\xi$. The prior distribution is also Gaussian, $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta)$. We assume unbiased measurements and priors centered on the origin, such that $\boldsymbol{\mu}_\xi = (0, 0)^\top = \boldsymbol{\mu}_\theta$, and parameterise the noise and prior covariances matrices by

$$\boldsymbol{\Sigma}_\xi = \begin{bmatrix} \sigma_{\xi_1}^2 & \rho_\xi \sigma_{\xi_1} \sigma_{\xi_2} \\ \rho_\xi \sigma_{\xi_2} \sigma_{\xi_1} & \sigma_{\xi_2}^2 \end{bmatrix},$$

$$\boldsymbol{\Sigma}_\theta = \begin{bmatrix} \sigma_{\theta_1}^2 & \rho_\theta \sigma_{\theta_1} \sigma_{\theta_2} \\ \rho_\theta \sigma_{\theta_2} \sigma_{\theta_1} & \sigma_{\theta_2}^2 \end{bmatrix},$$

where ρ_ξ and ρ_θ are the standard correlation coefficients in each case.

In this example, we take the opposite approach to that used in the univariate example, in that we retain the same likelihood function for each case considered and instead vary the form of the assumed prior, although in all cases the prior is centred on the origin of the parameter space. In particular, we assume throughout the true value $\boldsymbol{\theta}_* = (40, 40)^\top$, uncorrelated measurement noise with unit standard deviation, such that $\rho_\xi = 0$ and $\sigma_{\xi_1} = \sigma_{\xi_2} = 1$, and that the number of measurements is just $N = 1$, as this yields a circularly-symmetric bivariate Gaussian likelihood distribution of unit standard deviation, which is convenient for our investigations. We present results only for the BPR method, since the standard NS approach fails in all the cases considered.

Uncorrelated priors

We begin by considering uncorrelated priors, for which $\rho_\theta = 0$. In particular, we consider the two circularly-symmetric cases where $\{\sigma_{\theta_1}, \sigma_{\theta_2}\}$ equals to $\{4, 4\}$ and $\{2, 2\}$, and the intermediate non-circularly-symmetric case $\{2, 4\}$. These priors and the likelihood are plotted in Figure 6(a), which illustrates that all the priors are unrepresentative.

Figure 7 is a ‘corner plot’ showing the 1-dimensional and 2-dimensional marginal distributions of the joint ‘effective’ posterior on $(\theta_1, \theta_2, \beta)$ for each of the three priors described in Figure 6(a). In each case, the joint posterior again has the form of the product of two independent distributions on (θ_1, θ_2) and β , respectively, and is consistent with a marginal on β having the approximate form of a top-hat distribution in the range $[0, \beta_+]$. One also observes the expected evolution of this marginal across the three test cases, whereby β_+ gradually decreases as the likelihood is concentrated progressively further into the wings of the corresponding prior. Since $\beta_+ \lesssim 0.1$ in all three cases, one may confirm that each prior is indeed unrepresentative, and one may use the values of β_+ for each case to rank their severity.

Turning to the marginals on θ_1 and θ_2 in Figure 7, one sees that they are correctly centred on the mean of the corresponding true Gaussian posterior for each case. Moreover, the widths of the marginals are stable across the different priors and consistent with the width of the true Gaussian posterior in each case; the RMSE of the $\boldsymbol{\theta}_*$ estimate ≈ 0.05 in all cases. Thus, as in the univariate Gaussian likelihood example, the BPR method yields the correct inferences on (θ_1, θ_2) without any need for fine tuning.

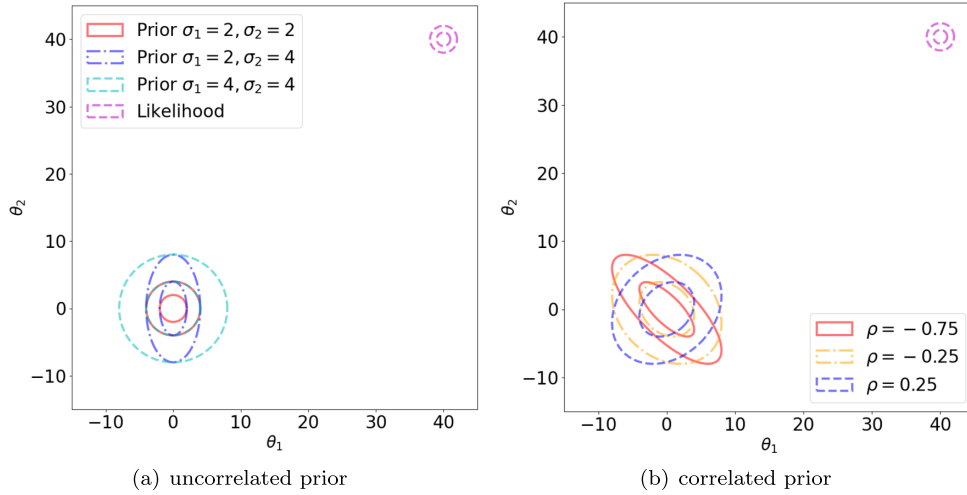


Figure 6: Illustration of the three test cases in the bivariate Gaussian likelihood example with uncorrelated/correlated priors. In the left panel, the cyan dashed, red solid and blue dot-dashed lines denote the Gaussian priors centred on the origin, with $\{\sigma_{\theta_1}, \sigma_{\theta_2}\}$ equals to $\{4, 4\}$, $\{2, 2\}$, and $\{2, 4\}$, respectively. The right panel shows correlated prior with $\{4, 4\}$ and correlation coefficients $\rho_{\theta} = \{-0.75, -0.25, 0.25\}$, respectively, denoted in a rainbow colour order. The pink dashed lines in the upper right corner denote the likelihood distribution. Each distribution contains two coloured contours corresponding to the 2σ (68%) and 3σ (95%) iso-probability levels.

One may also verify that the BPR method yields accurate evidence estimates in the above cases. Table 1 lists the mean and standard deviation of the estimated log-evidence over 10 realisations of the data for each of the three uncorrelated priors considered; it also lists the corresponding true evidences, estimated using standard quadrature techniques. One sees that in each case the log-evidence estimates obtained using BPR, after making the correction in (16), are all consistent with the true value, and have a standard deviation that remains stable for all the priors considered.

Correlated priors

We now consider priors with fixed standard deviations $\{\sigma_{\theta_1} = 4, \sigma_{\theta_2} = 4\}$, but three typical correlation coefficients $\rho_{\theta} = \{-0.75, -0.25, 0.25\}$. The three prior distributions are illustrated in Figure 6(b), together with the likelihood, from which one can see that all the priors are again unrepresentative.

Figure 8 is a ‘corner plot’ showing the 1-dimensional and 2-dimensional marginal distributions of the joint ‘effective’ posterior on $(\theta_1, \theta_2, \beta)$ for each of the three priors described, using the same colour order as in Figure 6(b). As in previous examples, the joint posterior in each case has the expected form of the product of two independent distributions on (θ_1, θ_2) and β , respectively, and the marginal on β is consistent with an

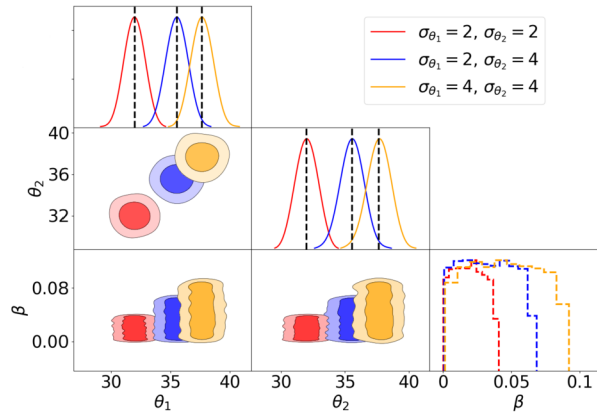


Figure 7: A ‘corner plot’ showing the 1-dimensional and 2-dimensional marginals of the joint posterior distribution on $(\theta_1, \theta_2, \beta)$ obtained by MultiNest using BPR for the bivariate Gaussian likelihood example with the likelihood and priors illustrated in Figure 6(a). The vertical dashed black lines indicate the mean of the true Gaussian posterior in each case.

Prior	True	BPR mean	BPR s.d.
$\{\sigma_{\theta_1} = 4, \sigma_{\theta_2} = 4\}$	-98.79	-98.70	0.88
$\{\sigma_{\theta_1} = 2, \sigma_{\theta_2} = 4\}$	-182.01	-182.01	1.24
$\{\sigma_{\theta_1} = 2, \sigma_{\theta_2} = 2\}$	-325.76	-325.66	1.49

Table 1: Comparison of the mean and standard deviation of the estimated log-evidence over 10 realisations of the data in the bivariate Gaussian likelihood example, obtained using the BPR method for a selection of uncorrelated priors. The ‘true’ value of the evidence in each case is also given, as estimated using standard quadrature techniques.

approximately top-hat distribution in the range $[0, \beta_+]$. The evolution of this marginal from $\rho_\theta = 0.25$ (blue) to $\rho_\theta = -0.75$ (red) is as expected, with β_+ gradually decreasing as the likelihood is concentrated further into the wings of the prior distribution, as the latter ‘rotates’ away from the peak of the likelihood. Moreover, since $\beta_+ \lesssim 0.15$ in every cases, one may conclude that all the priors are indeed unrepresentative and again use the β_+ values to rank their severity.

The marginals on θ_1 and θ_2 in Figure 8 are again correctly centred on the corresponding true Gaussian posterior for each case, the means of which are indicated by the vertical black dashed lines. The widths of the marginals are stable across the range of priors and consistent with the width of the true Gaussian posterior in each case; the RMSE of the θ_* estimate ≈ 0.06 in all cases. Hence the inferences on the parameters (θ_1, θ_2) using BPR are again accurate and robust in each case, without any tuning.

One may again verify that the BPR method yields accurate evidence estimates in the above cases. Table 2 lists the mean and standard deviation of the estimated log-evidence

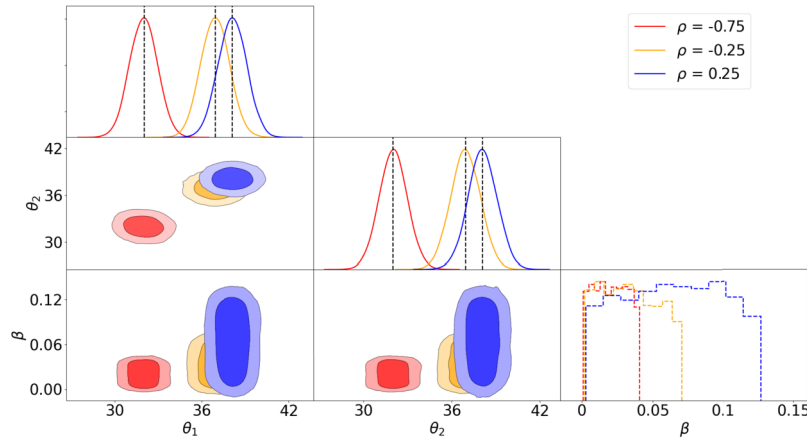


Figure 8: A ‘corner plot’ showing the 1-dimensional and 2-dimensional marginals of the joint posterior distribution on $(\theta_1, \theta_2, \beta)$ obtained by MultiNest using BPR for the bivariate Gaussian likelihood example with the likelihood and priors illustrated in Figure 6(b). The vertical dashed black lines indicate the mean of the true Gaussian posterior.

ρ_θ	True	BPR mean	BPR s.d.
-0.75	-324.39	-324.32	1.27
-0.25	-127.72	-127.65	0.85
0.25	-80.83	-80.75	0.98

Table 2: Comparison of the mean and standard deviation of the estimated log-evidence over 10 realisations of the data in the bivariate Gaussian likelihood example, obtained using the BPR method for a selection of priors with $\{\sigma_{\theta_1} = 4, \sigma_{\theta_2} = 4\}$ and correlation coefficient ρ_θ . The ‘true’ value of the evidence in each case is also given, as estimated using standard quadrature techniques.

values over 10 realisations of the data for each of the three correlated priors considered, and compares them with the corresponding true evidences estimated using standard quadrature techniques. As for the uncorrelated priors, the log-evidence estimates are all consistent with the true value, and have a stable standard deviation across all the priors considered.

5.3 Bivariate multi-modal likelihoods

We now consider bivariate multi-modal likelihoods consisting of an equal mixture of four identical but spatially separated Gaussians. Once again, for the sake of brevity, we consider only an uncorrelated Gaussian prior centred on the origin, with $\sigma_{\theta_1} = 4$ and $\sigma_{\theta_2} = 4$. Nonetheless, we do consider two different arrangements of the 4 modes of the likelihood relative to the centre of the prior. In particular we consider the *symmetric* and *asymmetric* arrangements illustrated in Figure 9, in which the circles denotes the 3-sigma iso-probability contour of each mode of the likelihood and the prior, respec-

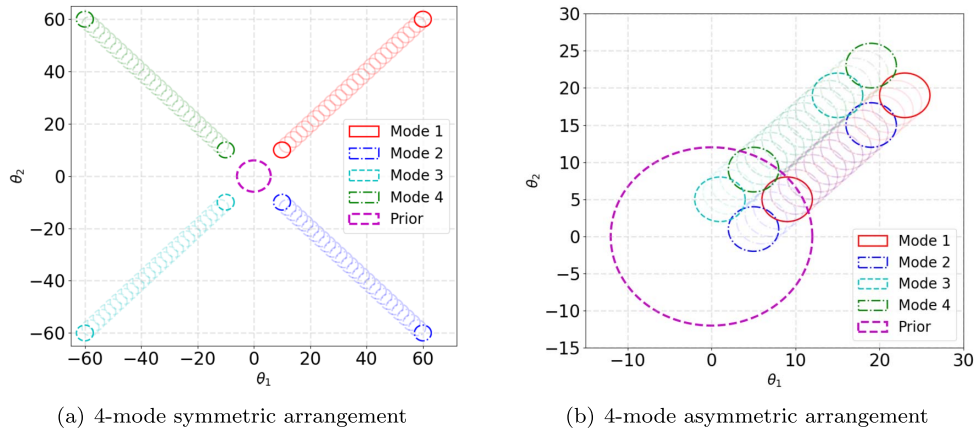


Figure 9: Illustration of the bivariate multimodal likelihoods considered, each of which consists of an equal mixture of four identical spatially-separated Gaussians for: (a) the symmetric arrangement; and (b) the asymmetric arrangement (relative to the prior). In each case, the dashed purple circle (centred at the origin) represents the 3-sigma iso-probability contour of the prior, and the other smaller circles represent the 3-sigma iso-probability contour of each of the four modes.

tively. As illustrated, for each arrangement, we consider a range of likelihoods for which the centres of the modes lie at varying distances from the origin. For the symmetric arrangement, each mode lies at the same distance from the origin, with positions ranging from $|\theta_1| = |\theta_2| = 5$ to $|\theta_1| = |\theta_2| = 60$. For the asymmetric arrangement, as shown in Figure 9(b), the centres of the 4 modes are placed a distance of 4 units from their geometric centre, the position of which ranges from $\theta_1 = \theta_2 = 5$ to 19. From Figure 9, one sees that for all the cases considered in the symmetric arrangement the prior is unrepresentative, whereas this is not true for some of the cases considered in the asymmetric arrangement. All other settings are identical to those used in the analysis of the bivariate Gaussian likelihood in Section 5.2.

Symmetric arrangement

Figure 10 summarises the performance of the BPR method in the symmetric arrangements illustrated in Figure 9(a). The upper row of panels (a)–(b) shows the effect of increasing the (equal) distance of the 4 mode centres from the origin for a fixed number of live points $N_{\text{live}} = 100$. The lower row of panels (c)–(d) shows the effect of increasing the number of live points N_{live} for fixed modes centres at $|\theta_1| = |\theta_2| = 40$.

In the left-hand column, i.e. panels (a) and (c), we plot $\log(\hat{\beta})$ (to accommodate the large dynamic range in the estimates of $\hat{\beta}$) obtained from the mean of the posterior samples associated with each of the four modes, separately. In panel (a), the value of $\log(\hat{\beta})$ is consistent across the four modes, as one might expect since they are symmet-

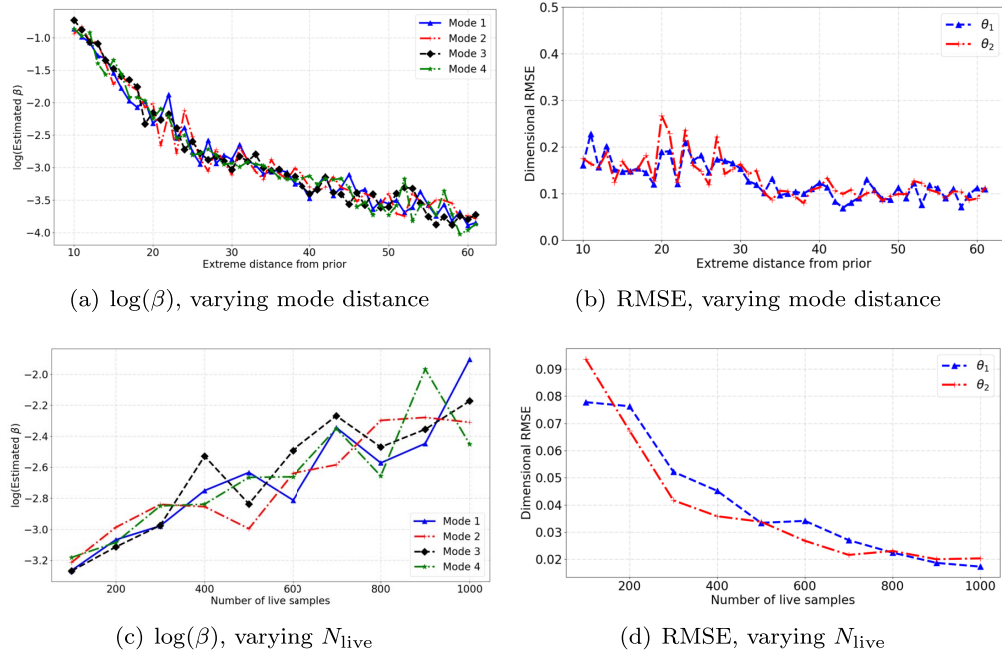


Figure 10: Performance of the BPR method applied to multimodal likelihoods consisting of an equal mixture of 4 identical but spatially-separated Gaussians in the *symmetric arrangement* illustrated in Figure 9(a). The upper row shows the effect of increasing the (equal) distance of the 4 mode centres from the origin for a fixed number of live points $N_{\text{live}} = 100$. The lower row shows the effect of increasing the number of live points N_{live} for fixed modes centres at $|\theta_1| = |\theta_2| = 40$. Panels (a) and (c) show $\log(\hat{\beta})$ obtained from posterior samples associated with each of the four modes, separately. Panels (b) and (d) show the RMSE in the estimation of θ_1 and θ_2 , averaged across the 4 modes.

rically arranged relative to the centre of the prior located at the origin. Moreover, as the (equal) distance of each mode from the origin increases, then the value of $\log(\hat{\beta})$ decreases monotonically (to within the statistical uncertainties of the nested sampling process). In panel (c), one again sees consistency on the values of $\hat{\beta}$ across the four modes, but that $\log(\hat{\beta})$ increases quasi-monotonically as N_{live} increases. This is expected since the failure of the nested sampling algorithm, and the consequent need for the posterior repartitioning characterised by low values of β , is reduced as N_{live} increases.

In the right-hand column of Figure 10, i.e. panels (b) and (d), we plot the RMSE in the estimation of θ_1 and θ_2 for each mode, but averaged across the 4 modes. In panel (b), one sees that the RMSE for θ_1 and θ_2 are consistent, as expected, and that the values remain stable in the range ~ 0.1 – 0.2 as the (equal) distance of the 4 modes from the origin varies. Panel (d) shows that the RMSE for θ_1 and θ_2 are again consistent and that these values decrease monotonically as N_{live} increases, as expected.

Asymmetric arrangement

Figure 11 summarises the performance of the BPR method in the asymmetric arrangements illustrated in Figure 9(b), and the quantities displayed in each panel are the same as those in Figure 10. The only differences here are that the panels (a) and (c) show the effects as a function of the distance of the geometric centre of the 4 modes from the origin, and that panels (b) and (d) show the RMSE for each mode separately. Further heatmap illustrations for this example can be found in Figure 15 in the Supplementary Material. This asymmetric modes scenario presents quite an extreme test of the BPR method, but equally illustrates some further key advantages over our original PR method presented in Chen et al. (2018).

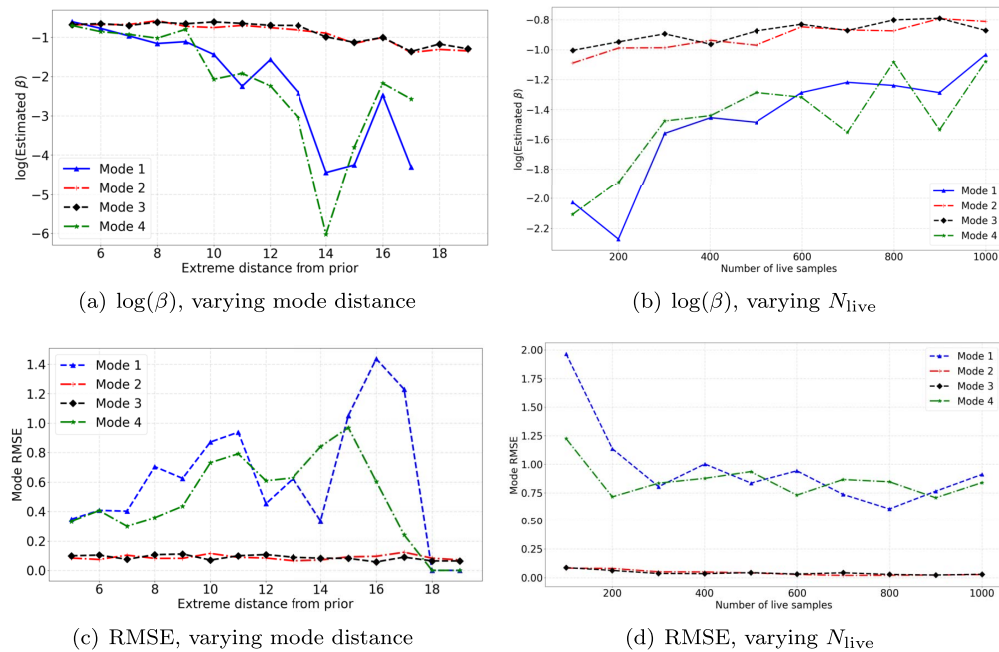


Figure 11: As for Figure 10, but for the *asymmetric arrangement* illustrated in Figure 9(b). In this case, the panels (a) and (c) show the effects as a function of the distance of the geometric centre of the 4 modes from the origin for $N_{\text{live}} = 100$, and the panels (b) and (d) show the effects of increasing N_{live} for fixed mode centres, of which the geometric centre is 15. Also, the panels (c) and (d) show the average RMSE in the estimation of θ_1 and θ_2 for each mode separately.

A unique feature of BPR in asymmetric multi-modal scenarios: In the upper row of panels (a) and (b), one sees that for each mode the behaviour of $\log(\hat{\beta})$ as a function of distance from the origin and N_{live} , respectively, is similar to that found for the symmetric arrangement, and for the same reasons as discussed above. By contrast, however, for the asymmetric arrangement the values of $\log(\hat{\beta})$ are not the same across

the 4 modes of the posterior. Instead, the $\log(\hat{\beta})$ values are very similar for modes 2 and 3, and broadly consistent for modes 1 and 4, with the values for the latter pair of modes being much smaller than those for the former pair. By inspecting Figure 9(b), one sees that this is to be expected, since modes 2 and 3 lie closer (but equidistant) from the centre of the prior at the origin, whereas modes 1 and 4 lie further away (but again equidistant). Thus, the prior is more unrepresentative for modes 1 and 4, which thus require a lower value of $\log(\hat{\beta})$.

This demonstrates an important feature of the BPR method, which is not shared by the original PR method, in that different regions of the posterior can be characterised by *different* ranges of β values, depending on how far into the wings of the prior the corresponding regions of likelihood lie, and this is accommodated in a fully automated manner. From panel (a), it is also worth noting that the difference between $\log(\hat{\beta})$ values for modes (2, 3) and modes (1, 4) is maximal for intermediate distances from the origin. This again makes sense as it corresponds to the point at which there is the largest difference between how representative the prior is for these two sets of modes, since the prior is changing most rapidly between them: for smaller distances (near the peak of the prior) or larger distances (in the wings of the prior), there is less change in the prior between the two sets of modes.

In panels (c) and (d) of Figure 11, one first sees from panel (c) that the average RMSE in the estimation of θ_1 and θ_2 , is low and insensitive to the distance from the origin for modes 2 and 3, which are closer the origin. By contrast, the RMSE is large and more volatile for modes 1 and 4, which lie further into the wings of the prior, and rises slowly with distance from the origin. One sees, however, that the RMSE drops to zero at very large distances; this occurs because, for the default setting of $N_{\text{live}} = 100$, one obtains no samples in modes 1 and 4. This therefore defines the limit of the applicability of the BPR method in this example, although the issue is resolved by increasing N_{live} . In panel (d), one sees that the RMSE is low and insensitive to N_{live} for modes 2 and 3, whereas modes 1 and 4 exhibit a larger and more volatile RMSE, albeit decreasing rapidly up to $N_{\text{live}} \approx 300$ and being relatively insensitive to N_{live} thereafter.

5.4 Higher-dimensional Gaussian likelihoods

Finally, we extend the bivariate Gaussian likelihood example in Section 5.2 to higher dimensions from 3D to 10D. We begin by adopting an analogous likelihood, corresponding to $N = 1$ data points and centered on $\theta_* = 40$ in all dimensions, but restrict our analysis to the case of a single circular-symmetric Gaussian prior, with $\sigma_\theta = 4$ in all dimensions. Once again, we consider only the BPR results as the standard NS method fails completely in this case. All results presented are based on 10 realisations of the data.

We first consider the constraints obtained on the parameters (θ, β) . Rather than presenting corner plots of the posterior in these higher dimensionalities, Figure 12(a)–(b) instead shows boxplots of the estimated value of β_+ (the 99% point of the unnormalised marginal posterior $\tilde{P}(\beta)$, as discussed in Section 4) and the RMSE of the estimate of θ_* (averaged across all dimensions), as a function of the dimensionality of the problem.

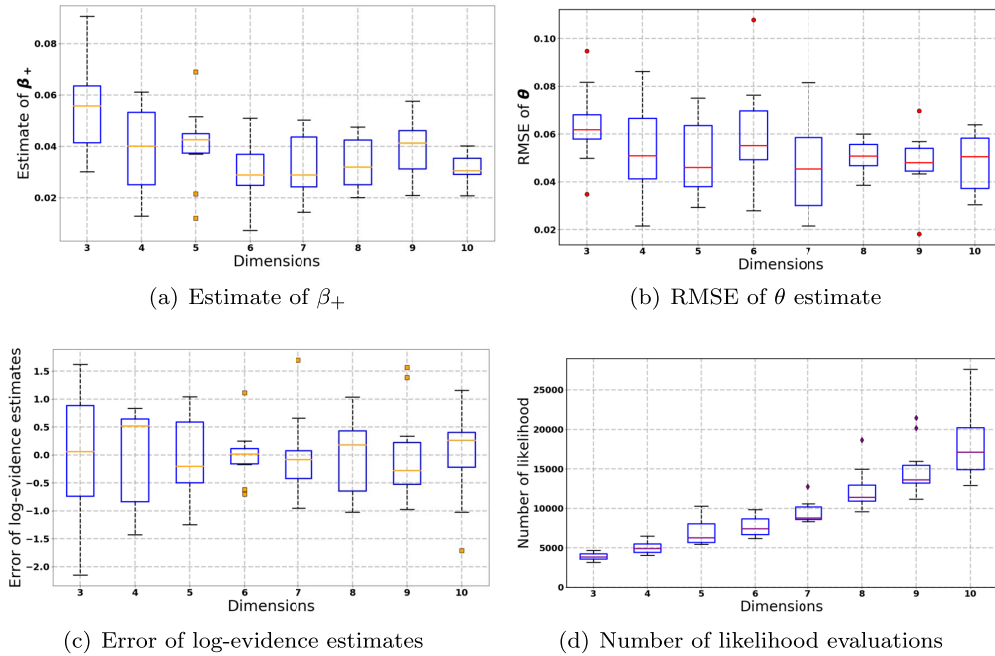


Figure 12: Boxplots for estimated parameter values using the BPR method applied to the higher-dimensional examples with $\theta_* = 40$, from 3D to 10D. The error bars denote the minimum and maximum values, whereas the boxes indicate the 25th to 75th quantiles. The orange/red line in the box represents the median value over 10 realisations of the data and the points denote outliers.

The estimated β_+ values are broadly consistent across different dimensionalities, with a mean of ≈ 0.05 and standard deviation ≈ 0.01 , with only a slight trend to smaller values as the dimensionality increases. This insensitivity to dimensionality is a result of the (effective) joint posterior having the product form (15). Similarly, the RMSE of the θ_* estimate is also seen to be broadly stable across the range of dimensionalities, with a mean value ≈ 0.05 , which is consistent with the accuracy obtained in the bivariate example, and standard deviation ≈ 0.015 .

Turning to the accuracy of the evidence estimates, Figure 12(c) shows a boxplot of the error in the log-evidence, over 10 realisations of the data, as a function of dimensionality. Once again, one sees that the distributions are relatively stable across different dimensionalities, with errors typically lying within one log-unit, which is consistent with the accuracies achieved in the bivariate example.

Finally, we consider the number of likelihood evaluations N_{like} required for MultiNest to converge. Figure 12(d) shows a boxplot of N_{like} over 10 realisations of the data, as a function of dimensionality. The mean value of N_{like} rises from ~ 3800 at 3D to ~ 18100 at 10D. This trend is consistent with that reported using the original PR method in

Chen et al. (2018), and follows roughly an $\mathcal{O}(n \log n)$ increase with dimensionality. The range of N_{like} values over the 10 realisations is also seen to widen slightly as the number of dimensions increases, but this is a relatively minor effect, at least up to 10D.

6 Conclusions

We have demonstrated that one may straightforwardly automate our previous prior repartitioning method for improving the robustness and efficiency of nested sampling in the presence of an unrepresentative prior. This is achieved by taking an explicitly Bayesian approach that treats the auxiliary parameter β in the power PR approach as a hyper-parameter that is estimated alongside the parameters of interest θ of the problem under consideration. Since this estimation process is performed within a single run of the nested sampling algorithm, the BPR method provides a substantial reduction in the computational requirements relative to the annealing schedule approach adopted in the original PR method. In addition to retaining all the advantages of the original scheme in providing reliable parameter constraints and evidence estimates, the BPR method adapts automatically to each problem and thus requires no tuning whatsoever. Indeed, by treating β as a hyper-parameter, one may use its resulting marginal to determine both the presence and the severity of an unrepresentative prior. We illustrate these properties in a range of numerical examples, from 1D to 10D, some exhibiting multi-modal likelihoods, with a selection of unrepresentative priors with varying degrees of severity. In particular, for multi-modal likelihoods, we show that different regions of the posterior may be characterised by different values of β , depending on how far into the wings of the prior the corresponding regions of likelihood lie. Moreover, this is accommodated in a fully automated manner and is an important feature of the BPR method that is not shared by our original PR method.

Perhaps most interestingly, we show that there is negligible computational overhead relative to standard nested sampling when the BPR method is used in cases where the prior is representative, and that it offers significant advantages in terms of efficiency and robustness even if the prior is only marginally unrepresentative. This suggests that the BPR should always be used in the application of nested sampling to Bayesian inference problems, irrespective of whether one suspects that the assumed priors may be unrepresentative for some data sets. Moreover, BPR might usefully be integrated directly into nested sampling algorithms and/or data analysis pipelines. It should be recalled, however, that the computational cost of BPR may be considerably increased if the normalising constant of the modified prior cannot be calculated analytically and so must be estimated numerically.

Supplementary Material

Supplementary materials for Bayesian posterior repartitioning for nested sampling (DOI: [10.1214/22-BA1323SUPP](https://doi.org/10.1214/22-BA1323SUPP); .pdf).

References

- Alsing, J. and Handley, W. (2021). “Nested sampling with any prior you like.” *arXiv e-prints*, [arXiv:2102.12478](https://arxiv.org/abs/2102.12478). 697, 701
- Baldock, R. J., Bernstein, N., Salerno, K. M., Pártay, L. B., and Csányi, G. (2017). “Constant-pressure nested sampling with atomistic dynamics.” *Physical Review E*, 96(4): 043311. 699
- Brewer, B. J., Pártay, L. B., and Csányi, G. (2011). “Diffusive nested sampling.” *Statistics and Computing*, 21(4): 649–656. [MR2826698](https://doi.org/10.1007/s11222-010-9198-8). doi: <https://doi.org/10.1007/s11222-010-9198-8>. 699
- Buchner, J. (2019). “Collaborative nested sampling: Big Data versus complex physical models.” *Publications of the Astronomical Society of the Pacific*, 131(1004): 108005. 699
- Buchner, J. (2021). “Nested Sampling Methods.” *arXiv preprint*, [arXiv:2101.09675](https://arxiv.org/abs/2101.09675). 699
- Chen, X., Feroz, F., and Hobson, M. (2022). “Supplementary materials for Bayesian posterior repartitioning for nested sampling.” *Bayesian Analysis*. doi: <https://doi.org/10.1214/22-BA1323SUPP>. 699
- Chen, X., Hobson, M., Das, S., and Gelderblom, P. (2018). “Improving the efficiency and robustness of nested sampling using posterior repartitioning.” *Statistics and Computing*, 1–16. [MR3955288](https://doi.org/10.1007/s11222-018-9841-3). doi: <https://doi.org/10.1007/s11222-018-9841-3>. 696, 697, 701, 702, 705, 706, 716, 719
- Chopin, N. and Robert, C. (2007). “Contemplating evidence: properties, extensions of, and alternatives to nested sampling.” Technical report, Technical Report 2007-46, CEREMADE, Université Paris Dauphine. 699
- Feroz, F. and Hobson, M. (2008). “Multimodal nested sampling: an efficient and robust alternative to Markov Chain Monte Carlo methods for astronomical data analyses.” *Monthly Notices of the Royal Astronomical Society*, 384(2): 449–463. 697
- Feroz, F., Hobson, M., and Bridges, M. (2009). “MultiNest: an efficient and robust Bayesian inference tool for cosmology and particle physics.” *Monthly Notices of the Royal Astronomical Society*, 398(4): 1601–1614. 696, 699, 701, 706, 707, 708
- Feroz, F., Hobson, M., Cameron, E., and Pettitt, A. (2013). “Importance nested sampling and the MultiNest algorithm.” *arXiv preprint*, [arXiv:1306.2144](https://arxiv.org/abs/1306.2144). 699
- Handley, W. (2019). “anesthetic: nested sampling visualisation.” *The Journal of Open Source Software*, 4(37). doi: <https://doi.org/10.21105/joss.01414>. 721
- Handley, W., Hobson, M., and Lasenby, A. (2015). “POLYCHORD: next-generation nested sampling.” *Monthly Notices of the Royal Astronomical Society*, 453(4): 4384–4398. 696, 697, 699
- Higson, E., Handley, W., Hobson, M., and Lasenby, A. (2018). “Dynamic nested sampling: an improved algorithm for parameter estimation and evidence calcu-

- lation.” *Statistics and Computing*. MR3994608. doi: <https://doi.org/10.1007/s11222-018-9844-0>. 699
- MacKay, D. (2003). *Information theory, inference and learning algorithms*. Cambridge university press. MR2012999. 699
- Skilling, J. (2006). “Nested Sampling for General Bayesian Computation.” *Bayesian Analysis*, 1(4): 833–860. MR2282208. doi: <https://doi.org/10.1214/06-BA127>. 695, 699
- Speagle, J. S. (2020). “dynesty: a dynamic nested sampling package for estimating Bayesian posteriors and evidences.” *Monthly Notices of the Royal Astronomical Society*, 493(3): 3132–3158. 699

Acknowledgments

The authors thank Will Handley and Lukas Hergt for their support with the powerful Python visualisation package `anesthetic` (Handley, 2019).