

Proximal nested sampling with data-driven priors for physical scientists

Jason D. McEwen^{1,2*}, Tobías I. Liaudat^{1,3}, Matthew A. Price¹, Xiaohao Cai⁴ and Marcelo Pereyra⁵

¹ Mullard Space Science Laboratory, University College London (UCL), Dorking, RH5 6NT, UK;

² Alan Turing Institute, London, NW1 2DB, UK;

³ Department of Computer Science, University College London (UCL), London, WC1E 6BT, UK;

⁴ School of Electronics and Computer Science, University of Southampton, Southampton, SO17 1BJ, UK;

⁵ School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, EH14 4AS, UK;

* Correspondence: jason.mcewen@gmail.com

Abstract: Proximal nested sampling was introduced recently to open up Bayesian model selection for high-dimensional problems such as computational imaging. The framework is suitable for models with a log-convex likelihood, which are ubiquitous in the imaging sciences. The purpose of this article is two-fold. First, we review proximal nested sampling in a pedagogical manner in an attempt to elucidate the framework for physical scientists. Second, we show how proximal nested sampling can be extended in an empirical Bayes setting to support data-driven priors, such as deep neural networks learned from training data.

Keywords: Bayesian model selection; nested sampling; proximal calculus.

1. Introduction

In much of the sciences not only is one interested in estimating the parameters of an underlying model, but deciding which model is best among a number of alternatives is of critical scientific interest. Bayesian model comparison provides a principled approach to model selection [1] that has found widespread application in the sciences [2].

Bayesian model comparison requires computation of the model evidence:

$$\mathcal{Z} = p(y|M) = \int dx p(y|x, M)p(x|M) = \int dx \mathcal{L}(x) \pi(x), \quad (1)$$

also called the marginal likelihood, where $y \in \mathbb{R}^m$ denotes data, $x \in \mathbb{R}^n$ parameters of interest, and M the model under consideration. We adopt the shorthand notation for the likelihood of $\mathcal{L}(x) = p(y|x, M)$ and prior of $\pi(x) = p(x|M)$. Evaluating the multi-dimensional integral of the model evidence is computationally challenging, particularly in high dimensions. While a number of highly successful approaches to computing the model evidence have been developed, such as nested sampling [e.g. 2–8] and the learned harmonic mean estimator [9–11], previous approaches do not scale to the very high-dimensional settings of computational imaging, which is our driving motivation.

The proximal nested sampling framework was introduced recently by a number of authors of the current article in order to open up Bayesian model selection for high-dimensional imaging problems [12]. Proximal nested sampling is suitable for models for which the likelihood is log-convex, which are ubiquitous in the imaging sciences. By restricting the class of models considered, it is possible to exploit structure of the problem to enable computation in very high-dimensional settings of $\mathcal{O}(10^6)$ and beyond.

Proximal nested sampling draws heavily on convex analysis and proximal calculus. In this article we present a pedagogical review of proximal nested sampling, sacrificing some mathematical rigor in an attempt to provide greater accessibility. We also provide a concise review of convexity and proximal calculus to introduce the background underpinning the framework. We assume the reader is familiar with nested sampling, hence we avoid repeating an introduction to nested sampling and instead refer the reader to other sources that provide excellent descriptions [2,3,8]. Finally, for the first time we show in an empirical



Citation: McEwen, J. D.; Liaudat, T.; Price, M. A.; Cai, X.; Pereyra, M. Proximal nested sampling with data-driven priors for physical scientists. *Preprints* 2023, 1, 0. <https://doi.org/>



Copyright: © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

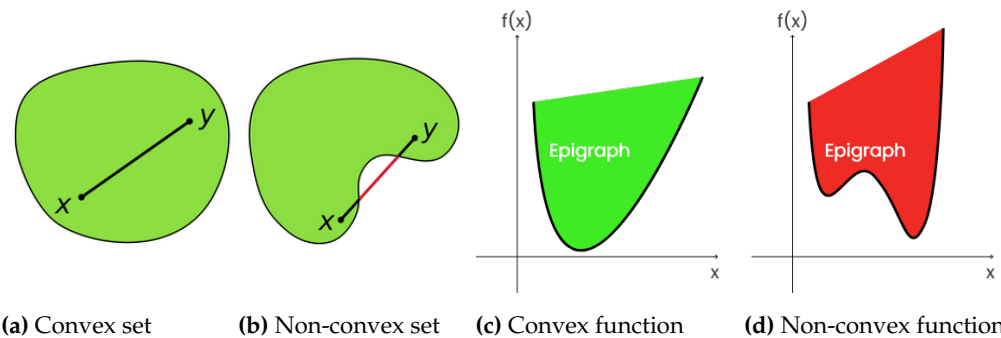
Bayes setting how proximal nested sampling can be extended to support data-driven priors, such as deep neural networks learned from training data.

2. Convexity and proximal calculus

We present a concise review of convexity and proximal calculus to introduce the background underpinning proximal nested sampling to make it more accessible.

2.1. Convexity

Proximal nested sampling draws on convexity, key concepts of which are illustrated in Figure 1. A set \mathcal{C} is convex if for any $x_1, x_2 \in \mathcal{C}$ and $\alpha \in (0, 1)$ we have $\alpha x_1 + (1 - \alpha)x_2 \in \mathcal{C}$. The epigraph of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by $\text{epi}(f) = \{(x, \gamma) \in \mathbb{R}^n \times \mathbb{R} \mid f(x) \leq \gamma\}$. The function f is convex if and only if its epigraph is convex. A convex function is lower semicontinuous if its epigraph is closed (*i.e.* includes its boundary).



(a) Convex set (b) Non-convex set (c) Convex function (d) Non-convex function

Figure 1. Proximal nested sampling considers likelihoods that are log-convex and lower semicontinuous. A lower semicontinuous convex function has a convex and closed epigraph.

2.2. Proximity operator

Proximal nested sampling leverages proximal calculus [13,14], a key component of which is the proximity operator, or prox. The proximity operator of the function f with parameter λ is defined by

$$\text{prox}_f^\lambda(x) = \arg \min_u [f(u) + \|u - x\|^2 / 2\lambda]. \quad (2)$$

The proximity operator maps a point x towards the minimum of f , while remaining in the proximity of the original point. The parameter λ controls how close the mapped point remains to x . An illustration is given in Figure 2.

The proximity operator can be considered as a generalisation of the projection onto a convex set. Indeed, the projection operator can be expressed as a prox by

$$\Pi_{\mathcal{C}}(x) = \arg \min_u [\chi_{\mathcal{C}}(u) + \|u - x\|^2 / 2], \quad (3)$$

with function f given by the characteristic function $\chi_{\mathcal{C}}(x) = \infty$ if $x \notin \mathcal{C}$ and zero otherwise.

2.3. Moreau-Yosida regularisation

The final component required in the development of proximal nested sampling is Moreau-Yosida regularisation [e.g. 14]. The Moreau-Yosida envelop of a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is given by the infimal convolution:

$$f^\lambda(x) = \inf_{u \in \mathbb{R}^n} f(u) + \frac{\|u - x\|^2}{2\lambda}. \quad (4)$$

The Moreau-Yosida envelope of a function can be interpreted as taking its convex conjugate, adding regularisation, before taking the conjugate again [14]. Consequently, it

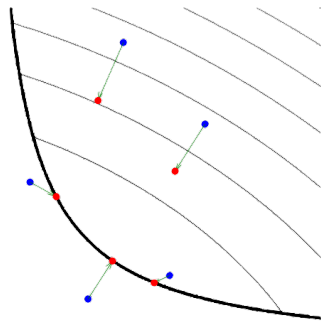


Figure 2. Illustration of the proximity operator (reproduced from [14]). The proximal operator maps the blue points to red points (*i.e.* from base to head of arrows). The thick black line defines the domain boundary, while the thin black lines define level-sets (iso-contours) of f . The proximity operator maps points towards the minimum of f , while remaining in the proximity of the original point.

provides a smooth regularised approximation of f , which is very useful to enable the use of gradient-based computational algorithms [*e.g.* 15].

The Moreau-Yosida envelope exhibits the following properties. First, λ controls the degree of regularisation with $f^\lambda(x) \rightarrow f(x)$ as $\lambda \rightarrow 0$. Second, the gradient of the Moreau-Yosida envelope of f can be computed through its prox by $\nabla f^\lambda(x) = (x - \text{prox}_f^\lambda(x))/\lambda$.

3. Proximal nested sampling

The challenge of nested sampling in high-dimensional settings is to sample from the prior distribution subject to a hard likelihood constraint [2,3,8]. Proximal nested sampling addresses this challenge for the case of log-convex likelihoods, which are widespread in computational imaging problems. In this section we review the proximal nested sampling framework [12] in a pedagogical manner, sacrificing some mathematical rigor in an attempt to improve readability and accessibility.

3.1. Constrained sampling formulation

Consider a prior and likelihood $\pi(x) \propto \exp(-f(x))$ and $\mathcal{L}(x) \propto \exp(-g(x))$, where the log-likelihood $g = -\log \mathcal{L}$ is a convex lower semicontinuous function. The log-prior $f = -\log \pi$ need only be differentiable or convex (it need not be convex if it is differentiable).

We consider sampling from the prior $\pi(x)$, such that $\mathcal{L}(x) > L^*$ for some likelihood value $L^* \geq 0$. Let $\iota_{L^*}(x)$ and $\chi_{L^*}(x)$ be the indicator function and characteristic function corresponding to this constraint, respectively, defined as

$$\iota_{L^*}(x) = \begin{cases} 1, & \mathcal{L}(x) > L^*, \\ 0, & \text{otherwise,} \end{cases} \quad \text{and} \quad \chi_{L^*}(x) = \begin{cases} 0, & \mathcal{L}(x) > L^*, \\ +\infty, & \text{otherwise.} \end{cases} \quad (5)$$

Since log is monotonic, $\mathcal{L}(x) > L^*$ is equivalent to $g(x) < \tau$ for $\tau = -\log L^*$. Explicitly define the convex set of the likelihood constraint by $\mathcal{B}_\tau = \{x \mid g(x) < \tau\}$. Then $\chi_{L^*}(x)$ is equivalent to $\chi_{\mathcal{B}_\tau}(x)$, where $\chi_{\mathcal{B}_\tau}(x) = \infty$ if $x \notin \mathcal{B}_\tau$ and zero otherwise.

Let $\pi_{L^*}(x) = \pi(x)\iota_{L^*}(x)$ represent the prior distribution with the hard likelihood constraint $\mathcal{L}(x) > L^*$. Since $\iota_{L^*}(x) = \exp(-\chi_{L^*}(x))$, then we have

$$-\log \pi_{L^*}(x) = -\log \pi(x) + \chi_{\mathcal{B}_\tau}(x). \quad (6)$$

To sample from the constrained prior we require sampling techniques that firstly can scale to high-dimensional settings and that secondly can support the convex constraint $\chi_{\mathcal{B}_\tau}(x)$.

3.2. Langevin MCMC sampling

Langevin Markov chain Monte Carlo (MCMC) sampling has been demonstrated to be highly effective at sampling in high-dimensional settings by exploiting gradient information

[15,16]. The Langevin stochastic differential equation associated with distribution $p(x)$ is a stochastic process defined by

$$dx(t) = \frac{1}{2} \nabla \log p(x(t)) dt + dw(t), \quad (7)$$

where $w(t)$ is Brownian motion. This process converges to $p(x)$ as time t increases and is therefore useful for generating samples from $p(x)$. In practice we compute a discrete-time approximation of $x(t)$ by the conventional Euler-Maruyama discretisation:

$$x^{(k+1)} = x^{(k)} + \frac{\delta}{2} \nabla \log p(x^{(k)}) + \sqrt{\delta} w^{(k+1)}, \quad (8)$$

where $w^{(k)}$ is a sequence of standard Gaussian random variables and δ is a step size.

Equation 8 provides a strategy for sampling in high-dimensions. However, notice that the updates rely on the score of the target distribution $\nabla \log p(\cdot)$. Nominally the target distribution must therefore be differentiable, which is not the case for our target of interest given by Equation 6. The prior may or may not be differentiable but the likelihood constraint certainly is not. Proximal versions of Langevin sampling have been developed to address the setting where the distribution is log-convex but not necessarily differentiable [15,16]. We follow a similar approach.

3.3. Proximal nested sampling framework

The proximal nested sampling framework follows by taking the constrained sampling formulation of Equation 6, adopting Langevin MCMC sampling of Equation 8, and applying Moreau-Yosida regularisation of Equation 4 to the convex constraint $\chi_{\mathcal{B}_\tau}(x)$ to yield a differentiable target. This strategy yields (see [12]) the update equation:

$$x^{(k+1)} = x^{(k)} + \frac{\delta}{2} \nabla \log \pi(x^{(k)}) - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})] + \sqrt{\delta} w^{(k+1)}, \quad (9)$$

where δ is the step size and λ is the Moreau-Yosida regularisation parameter.

Further intuition regarding proximal nested sampling can be gained by examining the term $v^{(k)} = -[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$, together with Figure 3. The vector $v^{(k)}$ points from the sample $x^{(k)}$ to its projection onto the likelihood constraint. If the sample $x^{(k)}$ is already in the likelihood-restricted prior support \mathcal{B}_τ , *i.e.* $x \in \mathcal{B}_\tau$, the term $v^{(k)}$ disappears and the Markov chain iteration simply involves the standard Langevin MCMC update. In contrast, if $x^{(k)}$ is not in \mathcal{B}_τ , *i.e.* $x \notin \mathcal{B}_\tau$, then a step is taken in the direction $v^{(k)}$, which acts to move the next iteration of the Markov chain in the direction of the projection of $x^{(k)}$ onto the convex set \mathcal{B}_τ . This term therefore acts to push the Markov chain back into the constraint set \mathcal{B}_τ if it wanders outside of it.¹

We have so far assumed that the (log) prior is differentiable (see Equation 9). This may not be the case, as is typical for sparsity-promoting priors (*e.g.* $-\log \pi(x) = \|\Psi^\top x\|_1 + \text{const.}$ for some wavelet dictionary Ψ). Then we make a Moreau-Yosida approximation of the log-prior, yielding the update equation:

$$x^{(k+1)} = x^{(k)} - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{-\log \pi}^\lambda(x^{(k)})] - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})] + \sqrt{\delta} w^{(k+1)}. \quad (10)$$

For notational simplicity here we have adopted the same regularisation parameter λ for each Moreau-Yosida approximation.

With the current formulation we are not guaranteed to recover samples from the prior subject to the hard likelihood constraint due to the approximation introduced in the Moreu-

¹ Note that proximal nested sampling has some similarity with Galilean [17] and constrained Hamiltonian [18] nested sampling. In these approaches Markov chains are also considered and if the Markov chain steps outside of the likelihood-constraint then it is reflected by an approximation of the shape of the boundary.

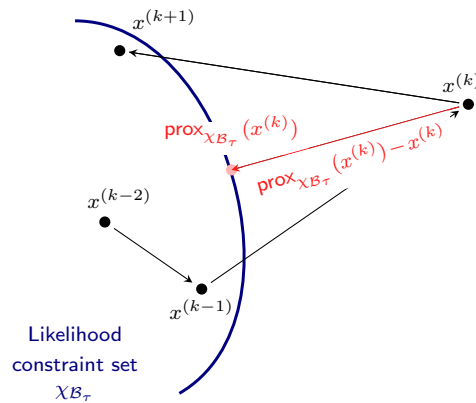


Figure 3. Diagram illustrating proximal nested sampling. If a sample $x^{(k)}$ outside of the likelihood constraint is considered, then proximal nested sampling introduces a term in the direction of the projection of $x^{(k)}$ onto the convex set defining the likelihood constraint, thereby acting to push the Markov chain back into the constraint set \mathcal{B}_τ if it wanders outside of it. A subsequent Metropolis-Hastings step can be introduced to enforce strict adherence to the convex likelihood constraint.

Yosida regularisation and due to the approximation in discretising the underlying Langevin stochastic differential equation. We therefore introduce a Metropolis-Hastings correction step to eliminate the bias introduced by these approximations and ensure convergence to the required target density (see [12] for further details).

Finally, we adopt this strategy for sampling from the constrained prior in the standard nested sampling strategy to recover the proximal nested sampling framework. The algorithm can be initialised with samples from the prior as described by the update equations above but with the likelihood term removed, *i.e.* with $[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})] \rightarrow 0$.

3.4. Explicit forms of proximal nested sampling

While we have discussed the general framework for proximal nested sampling, we have yet to address the issue of computing the proximity operators involved. As Equation 2 demonstrates, computing proximity operators involves solving an optimisation problem. Only in certain cases are closed form solutions available [13]. Explicit forms of proximal nested sampling must therefore be considered for the problem at hand.

We focus on a common high-dimensional inverse imaging problem where we acquire noisy observations $y = \Phi x + n$, of an underlying image x via some measurement model Φ , in the presence of Gaussian noise n (without loss of generality we consider independent and identically distributed noise here). We consider a Gaussian negative likelihood, $-\log \mathcal{L}(x) = \|y - \Phi x\|_2^2 / 2\sigma^2 + \text{const.}$, and a sparsity-promoting prior, $-\log \pi(x) = \mu \|\Psi^\dagger x\|_1 + \text{const.}$, for some wavelet dictionary Ψ . The prox of the prior can be computed in closed-form by [13]

$$\text{prox}_{-\log \pi}^\lambda(x) = x + \Psi(\text{soft}_{\lambda\mu}(\Psi^\dagger x) - \Psi^\dagger x), \quad (11)$$

where $\text{soft}_\lambda(\cdot)$ is the soft thresholding function with threshold λ (recall μ is the scale of the sparsity-promoting prior, *i.e.* the regularisation parameter, defined above). However, the prox of the likelihood is not so straightforward. The prox for the likelihood can be recast as a saddle-point problem that can be solved iteratively by a primal dual method initialised by the current sample position (see [12] for further details):

$$1. \quad z^{(i+1)} = z^{(i)} + \delta_1 \Phi \bar{x}^{(i)} - \text{prox}_{\chi_{\mathcal{B}'_\tau}}(z^{(i)} + \delta_1 \Phi \bar{x}^{(i)}),$$

$$\text{where } \text{prox}_{\chi_{\mathcal{B}'_\tau}}(z) = \text{proj}_{\mathcal{B}'_\tau}(z) = \begin{cases} z, & \text{if } z \in \mathcal{B}'_\tau, \\ \frac{z-y}{\|z-y\|_2} \sqrt{2\tau\sigma^2} + y, & \text{otherwise;} \end{cases}$$

2. $x^{(i+1)} = (x' + x^{(i)} - \delta_2 \Phi^\dagger z^{(i+1)})/2$;
3. $\bar{x}^{(i+1)} = x^{(i+1)} + \delta_3 (x^{(i+1)} - x^{(i)})$.

Combining these algorithms to efficiently compute prox operators with the proximal nested sampling framework, we can compute the model evidence to perform Bayesian model comparison in high-dimensional settings. We can also obtain posterior distributions with the usual weighted samples from the dead points of nested sampling. This allows one to recover, for example, point estimates such as the posterior mean image.

4. Deep data-driven priors

While hand-crafted priors, such as wavelet-based sparsity promoting priors, are common in computational imaging, they provide only limited expressivity. If example images are available an empirical Bayes approach with data-driven priors can be taken, where the prior is learned from training data. Since proximal nested sampling requires only the log-likelihood to be convex, complex data-driven priors, such as represented by deep neural networks, can be integrated into the framework. Through Tweedie's formula we describe how proximal nested sampling can be adapted to support data-driven priors, opening up Bayesian model selection for data-driven approaches. We take a similar approach to [19], where data-driven priors are integrated into Langevin MCMC sampling strategies, although in that work model selection is not considered.

4.1. Tweedie's formula and data-driven priors

Tweedie's formula is a remarkable result in Bayesian estimation credited to personal correspondence with Maurice Kenneth Tweedie [20]. Tweedie's formula has gained renewed interest in recent years [19,21–23] due to its connection to score matching [24–26] and denoising diffusion models [27,28], which as of this writing provide state-of-the-art performance in deep generative modelling.

Tweedie's result follows by considering the following scenario. Consider x sampled from a prior distribution $q(\cdot)$ and noisy observations $z \sim \mathcal{N}(x, \sigma^2 I)$. Tweedie's formula gives the posterior expectation of x given z as

$$E(x|z) = z + \sigma^2 \nabla \log p(z), \quad (12)$$

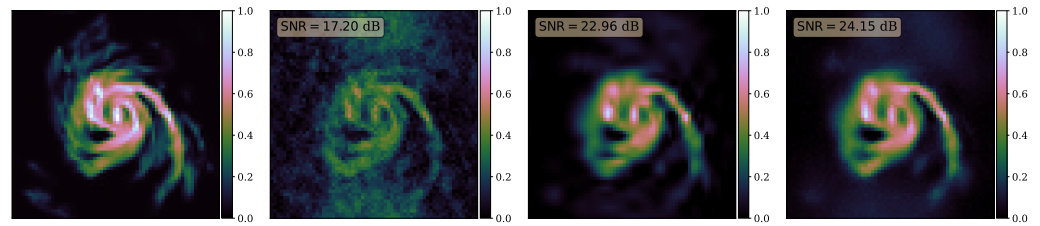
where $p(z)$ is the marginal distribution of z (for further details see, *e.g.*, [21]). The critical advantage of Tweedie's formula is that it does not require knowledge of the underlying distribution $q(\cdot)$ but rather only the marginalised distribution of the observation. Equation 12 can be interpreted as a denoising strategy to estimate x from noisy observations z . Moreover, Tweedie's formula can also be used to relate a denoiser (potentially a trained deep neural network) to the score $\nabla \log p(z)$.

In a data-driven setting, where the underlying prior is implicitly specified by training data (which are considered to be samples from the prior), there is no guarantee that the underlying prior, and therefore the posterior, is well-suited for gradient-based Bayesian computation such as Langevin sampling, *e.g.* it may not be differentiable. Therefore we consider a regularised version of the prior defined by Gaussian smoothing:

$$p_\epsilon(x) = (2\pi\epsilon)^{-n/2} \int dx' \exp(-\|x - x'\|_2^2 / (2\epsilon)) q(x'). \quad (13)$$

This regularisation can also be viewed as adding a small amount of regularising Gaussian noise. We can therefore leverage Tweedie's formula to relate the regularised prior distribution $p_\epsilon(x)$ to a denoiser D_ϵ trained to recover x from noisy observations $x_\epsilon \sim \mathcal{N}(x, \epsilon I)$, *i.e.* the score of the regularised prior can be related to the denoiser by

$$\nabla \log p_\epsilon(x) = \epsilon^{-1} (D_\epsilon(x) - x). \quad (14)$$



(a) Ground truth (b) Dirty (c) Hand-crafted prior (d) Data-driven prior

Figure 4. Results of radio interferometric imaging reconstruction problem. (a) Ground truth galaxy image. (b) Dirty reconstruction based on pseudo-inverting the measurement operator Φ . (c) Posterior mean reconstruction computed from proximal nested samples for the hand-crafted wavelet-sparsity prior. (d) Posterior mean reconstruction for the data-driven prior based on a deep neural network (DnCNN) trained on example images. Reconstruction SNR is shown on each image. The computed SNR levels demonstrate that the data-driven prior results in a superior reconstruction quality, although this may not be obvious from a visual assessment of the reconstructed images. Computing the reconstructed SNR requires knowledge of the ground truth, which is not available in realistic settings. The Bayesian model evidence proves a way to compare the hand-crafted and data-driven models without requiring knowledge of the ground truth. For this example the Bayesian evidence correctly selects the data-driven prior as the best model.

Denoisers are commonly integrated in proximal optimisation algorithms in replace of proximity operators, giving rise to so-called plug-and-play (PnP) methods [29,30] and, more recently, also into Bayesian computational algorithms [19]. Typically denoisers are represented by deep neural networks, which can be trained by injecting a small amount of noise in training data and learning to denoise the corrupted data. While a noise level ϵ needs to be chosen, as discussed above this is considered a regularisation of the prior and so the denoiser need not be trained on the noise level of a problem at hand. In this manner, the same denoiser can be used for multiple subsequent problems (hence the PnP name). The learned score of the regularised prior inherits the same properties as the denoiser, such as smoothness, hence the denoiser should be considered carefully. Well-behaved denoisers have been considered already in PnP methods (in order to provide convergence guarantees) and a popular approach for imaging problems is the DnCNN model [30], which is based on a deep convolutional neural network, and that is (Lipschitz) continuous.

4.2. Proximal nested sampling with data-driven priors

By Tweedie's formula the standard proximal nested sampling update of Equation 9 can be revised to integrate a learned denoiser, yielding

$$x^{(k+1)} = x^{(k)} - \frac{\alpha\delta}{2\epsilon} [x - D_\epsilon(x^{(k)})] - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{\chi_{B_\tau}}(x^{(k)})] + \sqrt{\delta}w^{(k+1)}, \quad (15)$$

where we have included a regularisation parameter α that allows us to balance the influence of the prior and the data fidelity terms [19]. We typically consider a deep convolutional neural network based on the DnCNN model [30] since it is (Lipschitz) continuous and has been demonstrated to perform very well in PnP settings [19,30]. Again, this sampling strategy can then be integrated into the standard nested sampling framework.

We can therefore support data-driven priors in the proximal nested sampling framework by integrating a deep denoiser that learns to denoise training data, using Tweedie's formula to relate this to the score of a regularised data-driven prior.

5. Numerical experiments

The new methodology presented allows us to perform Bayesian model comparison between a data-driven and hand-crafted prior(validation of proximal nested sampling in a setting where the ground truth can be computed directly has been performed already

[12]). We consider a simple radio interferometric imaging reconstruction problem as an illustration. We assume the same observational model as Section 3.4, with white Gaussian noise giving a signal-to-noise ratio (SNR) of 15dB. The measurement operator Φ is a masked Fourier transform as a simple model of a radio interferometric telescope. The mask is built by randomly selecting 50% of the Fourier coefficients. A Gaussian likelihood is used in both models. For the hand-crafted prior we consider a sparsity-promoting prior using a Daubechis 6 wavelet dictionary. We base the data-driven prior on a DnCNN [30] model trained on galaxy images extracted from the IllustrisTNG simulations [31]. We also consider an IllustrisTNG galaxy simulation, not used in training, as the ground truth test image. We generate samples following the proximal nested sampling strategies of Equation 10 and Equation 15 for the hand-crafted and data-driven priors, respectively. Posterior inferences (e.g. posterior mean image) and the model evidence can then be computed from nested sampling samples in the usual manner. The step size δ is set to 10^{-7} , the Moreau-Yosida regularisation parameter λ to 5×10^{-7} , and the regularisation strength of wavelet-based model μ to 5×10^4 . We consider noise level $\epsilon \simeq 8.34$ and set the regularisation parameter α of the data-driven prior to 3.5×10^{-7} . For the nested sampling methods, the number of live and dead samples is set to 10^2 and 2.5×10^3 , respectively. For the Langevin sampling, we use a thinning factor of 20 and set the number of burn-in iterations to 10^2 .

Results are presented in Figure 4. The data-driven prior results in a superior reconstruction with an improvement in SNR of 1.2dB, although it may be difficult to tell simply from visual inspection of the recovered images. Computing the SNR of the reconstructed images requires knowledge of the ground truth, which clearly is not accessible in realistic settings involving real observational data. The Bayesian model evidence, computed by proximal nested sampling, proves a way to compare the hand-crafted and data-driven models without requiring knowledge of the ground truth and is therefore applicable in realistic scenarios. We compute log evidences of -2.96×10^3 for the hand-crafted prior and -1.35×10^3 for the data-driven prior. Consequently, the data-driven model is preferred by the model evidence, which agrees with the SNR levels computed from the ground truth. These results are all as one might expect since learned data-driven priors are more expressive than hand-crafted priors and can better adapt to model high-dimensional images.

6. Conclusions

Proximal nested sampling leverages proximal calculus to extend nested sampling to high-dimensional settings for problems involving log-convex likelihoods, which are ubiquitous in computational imaging. The purpose of this article is two-fold. First, we review proximal nested sampling in a pedagogical manner in an attempt to elucidate the framework for physical scientists. Second, we show how proximal nested sampling can be extended in an empirical Bayes setting to support data-driven priors, such as deep neural networks learned from training data. We show only preliminary results for learned proximal nested sampling and will present a more thorough study in a follow-up article.

Author Contributions: Conceptualization, J.D.M. and M.P.; methodology, J.D.M., X.C. and M.P.; software, T.I.L., M.A.P. and X.C.; validation, T.I.L., M.A.P. and X.C.; resources, J.D.M.; data curation, M.A.P.; writing—original draft preparation, J.D.M. and T.I.L.; writing—review and editing, J.D.M., T.I.L., M.A.P., X.C. and M.P.; supervision, J.D.M.; project administration, J.D.M.; funding acquisition, J.D.M., M.A.P. and M.P. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by EPSRC grant number EP/W007673/1.

Data Availability Statement: The ProxNest code and experiments are available at <https://github.com/astro-informatics/proxnest>.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Robert, C.P. *The Bayesian Choice*; Springer-Verlag New York, 2007.

2. Ashton, G.; Bernstein, N.; Buchner, J.; Chen, X.; Csányi, G.; Fowlie, A.; Feroz, F.; Griffiths, M.; Handley, W.; Habeck, M.; et al. Nested sampling for physical scientists. *Nature Reviews Methods Primers* **2022**, *2*, 39.
3. Skilling, J. Nested sampling for general Bayesian computation. *Bayesian Analysis* **2006**, *1*, 833–859.
4. Mukherjee, P.; Parkinson, D.; Liddle, A.R. A nested sampling algorithm for cosmological model selection. *The Astrophysical Journal* **2006**, *638*, L51–L54.
5. Feroz, F.; Hobson, M.P. Multimodal nested sampling: an efficient and robust alternative to MCMC methods for astronomical data analysis. *Monthly Notices of the Royal Astronomical Society (MNRAS)* **2008**, *384*, 449–463.
6. Feroz, F.; Hobson, M.P.; Bridges, M. MULTINEST: an efficient and robust Bayesian inference tool for cosmology and particle physics. *Monthly Notices of the Royal Astronomical Society (MNRAS)* **2009**, *398*, 1601–1614.
7. Handley, W.J.; Hobson, M.P.; Lasenby, A.N. POLYCHORD: nested sampling for cosmology. *Monthly Notices of the Royal Astronomical Society: Letters* **2015**, *450*, L61–L65.
8. Buchner, J. Nested sampling methods. *arXiv preprint arXiv:2101.09675* **2021**, [[arXiv:2101.09675](https://arxiv.org/abs/2101.09675)].
9. McEwen, J.D.; Wallis, C.G.R.; Price, M.A.; Docherty, M.M. Machine learning assisted Bayesian model comparison: the learnt harmonic mean estimator. *Statistics & Computing*, submitted, **2022**, [[arXiv:2111.12720](https://arxiv.org/abs/2111.12720)].
10. Spurio Mancini, A.; Docherty, M.M.; Price, M.A.; McEwen, J.D. Bayesian model comparison for simulation-based inference. *RASTI*, submitted **2022**, [[arXiv:2207.04037](https://arxiv.org/abs/2207.04037)].
11. Polanska, A.; Price, M.A.; Spurio Mancini, A.; McEwen, J.D. Learned harmonic mean estimation of the marginal likelihood with normalising flows. *MaxEnt*, submitted **2023**, [[arXiv:2307.00048](https://arxiv.org/abs/2307.00048)].
12. Cai, X.; McEwen, J.D.; Pereyra, M. Proximal nested sampling for high-dimensional Bayesian model selection. *Statistics & Computing* **2022**, *32*, [[arXiv:2106.03646](https://arxiv.org/abs/2106.03646)].
13. Combettes, P.; Pesquet, J.C. *Proximal splitting methods in signal processing*; Springer: New York, 2011; pp. 185–212.
14. Parikh, N.; Boyd, S. Proximal algorithms. *Foundations and Trends in Optimization* **2013**, *1*, 123–231.
15. Pereyra, M. Proximal Markov chain Monte Carlo algorithms. *Statistics and Computing* **2016**, *26*, 745–760.
16. Durmus, A.; Moulines, E.; Pereyra, M. Efficient Bayesian computation by proximal Markov chain Monte Carlo: when Langevin meets Moreau. *SIAM Journal Imaging Sciences* **2018**, *1*, 473–506.
17. Skilling, J. Bayesian computation in big spaces-nested sampling and Galilean Monte Carlo. In Proceedings of the AIP Conference Proceedings 31st. American Institute of Physics, 2012, Vol. 1443, pp. 145–156.
18. Betancourt, M. Nested sampling with constrained hamiltonian monte carlo. In Proceedings of the AIP Conference Proceedings. American Institute of Physics, 2011, Vol. 1305, pp. 165–172.
19. Laumont, R.; Bortoli, V.D.; Almansa, A.; Delon, J.; Durmus, A.; Pereyra, M. Bayesian imaging using Plug & Play priors: when Langevin meets Tweedie. *SIAM Journal on Imaging Sciences* **2022**, *15*, 701–737.
20. Robbins, H. An Empirical Bayes Approach to Statistics. In Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, 1956, Astronomical Society of the Pacific Conference Series, pp. 157–163.
21. Efron, B. Tweedie’s formula and selection bias. *Journal of the American Statistical Association* **2011**, *106*, 1602–1614.
22. Kim, K.; Ye, J.C. Noise2score: tweedie’s approach to self-supervised image denoising without clean images. *Advances in Neural Information Processing Systems* **2021**, *34*, 864–874.
23. Chung, H.; Sim, B.; Ryu, D.; Ye, J.C. Improving diffusion models for inverse problems using manifold constraints. *arXiv preprint arXiv:2206.00941* **2022**.
24. Sohl-Dickstein, J.; Weiss, E.; Maheswaranathan, N.; Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In Proceedings of the International Conference on Machine Learning. PMLR, 2015, pp. 2256–2265.
25. Song, Y.; Ermon, S. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems* **2019**, *32*.
26. Song, Y.; Ermon, S. Improved techniques for training score-based generative models. *Advances in neural information processing systems* **2020**, *33*, 12438–12448.
27. Song, Y.; Sohl-Dickstein, J.; Kingma, D.P.; Kumar, A.; Ermon, S.; Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456* **2020**.
28. Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; Ommer, B. High-resolution image synthesis with latent diffusion models. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 10684–10695.
29. Venkatakrisnan, S.V.; Bouman, C.A.; Wohlberg, B. Plug-and-play priors for model based reconstruction. In Proceedings of the 2013 IEEE Global Conference on Signal and Information Processing. IEEE, 2013, pp. 945–948.
30. Ryu, E.; Liu, J.; Wang, S.; Chen, X.; Wang, Z.; Yin, W. Plug-and-play methods provably converge with properly trained denoisers. In Proceedings of the International Conference on Machine Learning. PMLR, 2019, pp. 5546–5557.
31. Nelson, D.; Springel, V.; Pillepich, A.; Rodriguez-Gomez, V.; Torrey, P.; Genel, S.; Vogelsberger, M.; Pakmor, R.; Marinacci, F.; Weinberger, R.; et al. The IllustrisTNG Simulations: Public Data Release. *Computational Astrophysics and Cosmology* **2019**, *6*, 2.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.