

Bayesian model selection

in cosmology and astrophysics

Jason D. McEwen

www.jasonmcewen.org

Mullard Space Science Laboratory (MSSL), University College London (UCL)

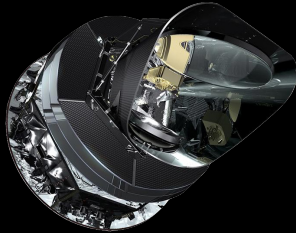
3rd IMA Conference on Inverse Problems from Theory to Application, May 2022

Observable Universe

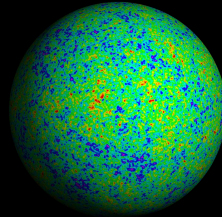


Cosmic microwave background (CMB) radiation

What is the origin of structure in our Universe?



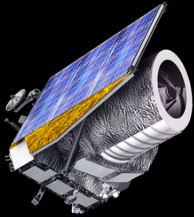
Planck satellite



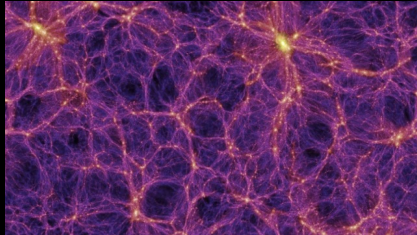
CMB

Large-scale structure of the Universe

What is the nature of dark energy?



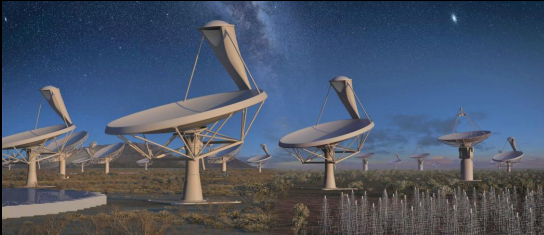
Euclid satellite



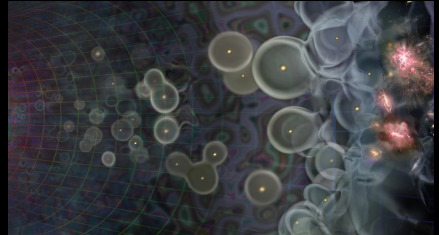
Large-scale structure

Epoch of reionisation

How did the first luminous objects in the Universe form?



Square Kilometre Array (SKA)



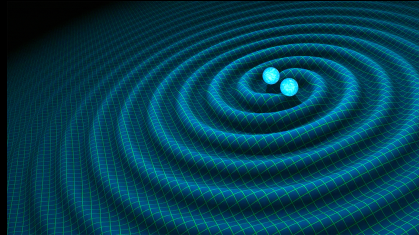
Ionised bubbles in neutral hydrogen

Gravitational waves

What are the physical processes responsible for an observed gravitational wave signal?



LIGO



Merging black holes

Model selection in cosmology and astrophysics

These are questions of **model selection**.

Model selection in cosmology and astrophysics

These are questions of **model selection**.

In cosmology we **cannot perform experiments** but just have **one Universe** to observe.

Model selection in cosmology and astrophysics

These are questions of **model selection**.

In cosmology we **cannot perform experiments** but just have **one Universe** to observe.

In astrophysics, we again cannot perform experiments, but may have a small number of observations of similar processes.

Model selection in cosmology and astrophysics

These are questions of **model selection**.

In cosmology we **cannot perform experiments** but just have **one Universe** to observe.

In astrophysics, we again cannot perform experiments, but may have a small number of observations of similar processes.

⇒ **Bayesian model selection**

1. Bayesian model selection
2. Learnt harmonic mean estimator for likelihood-based model selection
3. Learnt harmonic mean estimator for simulation-based model selection
4. Proximal nested sampling for high-dimensional model selection

Bayesian model selection

Bayesian inference: parameter estimation

Bayes' theorem

$$p(\theta | y, M) = \frac{p(y | \theta, M) p(\theta | M)}{p(y | M)}$$

Diagram illustrating Bayes' theorem with components labeled:

- Posterior: $p(\theta | y, M)$ (yellow box)
- Likelihood: $p(y | \theta, M)$ (orange box)
- Prior: $p(\theta | M)$ (blue box)
- Evidence: $p(y | M)$ (red box)

for parameters θ , model M and observed data y .

Bayesian inference: parameter estimation

Bayes' theorem

$$p(\theta | y, M) = \frac{p(y | \theta, M) p(\theta | M)}{p(y | M)} = \frac{\mathcal{L}(\theta) \pi(\theta)}{z}$$

The diagram illustrates Bayes' theorem with colored boxes and labels:

- The posterior $p(\theta | y, M)$ is shown in a yellow box.
- The likelihood $p(y | \theta, M)$ is shown in an orange box.
- The prior $p(\theta | M)$ is shown in a blue box.
- The evidence $p(y | M)$ is shown in a red box.
- The likelihood $\mathcal{L}(\theta)$ is shown in an orange box.
- The prior $\pi(\theta)$ is shown in a blue box.
- The evidence z is shown in a red box.

for parameters θ , model M and observed data y .

Bayesian inference: parameter estimation

Bayes' theorem

$$p(\theta | y, M) = \frac{p(y | \theta, M) p(\theta | M)}{p(y | M)} = \frac{\mathcal{L}(\theta) \pi(\theta)}{z}$$

The diagram illustrates Bayes' theorem with colored boxes and labels:

- The posterior $p(\theta | y, M)$ is shown in a yellow box.
- The likelihood $p(y | \theta, M)$ is in an orange box, labeled "likelihood".
- The prior $p(\theta | M)$ is in a blue box, labeled "prior".
- The evidence $p(y | M)$ is in a red box, labeled "evidence".
- The likelihood $\mathcal{L}(\theta)$ is in an orange box, labeled "likelihood".
- The prior $\pi(\theta)$ is in a blue box, labeled "prior".
- The evidence z is in a red box, labeled "evidence".

for parameters θ , model M and observed data y .

For **parameter estimation**, typically draw samples from the posterior by *Markov chain Monte Carlo (MCMC)* sampling.

Bayesian inference: model selection

By Bayes' theorem for model M_j :

$$p(M_j | y) = \frac{p(y | M_j)p(M_j)}{\sum_j p(y | M_j)p(M_j)} .$$

Bayesian inference: model selection

By Bayes' theorem for model M_j :

$$p(M_j | y) = \frac{p(y | M_j)p(M_j)}{\sum_j p(y | M_j)p(M_j)}.$$

For **model selection**, consider posterior model odds:

$$\boxed{\frac{p(M_1 | y)}{p(M_2 | y)}} = \boxed{\frac{p(y | M_1)}{p(y | M_2)}} \times \boxed{\frac{p(M_1)}{p(M_2)}}.$$

posterior odds Bayes factor prior odds

Bayesian inference: model selection

By Bayes' theorem for model M_j :

$$p(M_j | y) = \frac{p(y | M_j)p(M_j)}{\sum_j p(y | M_j)p(M_j)}.$$

For **model selection**, consider posterior model odds:

$$\boxed{\frac{p(M_1 | y)}{p(M_2 | y)}} = \boxed{\frac{p(y | M_1)}{p(y | M_2)}} \times \boxed{\frac{p(M_1)}{p(M_2)}}.$$

posterior odds Bayes factor prior odds

Must compute the **Bayesian model evidence** or **marginal likelihood** given by the normalising constant

$$z = p(y | M) = \int d\theta \mathcal{L}(\theta) \pi(\theta).$$

Bayesian inference: model selection

By Bayes' theorem for model M_j :

$$p(M_j | y) = \frac{p(y | M_j)p(M_j)}{\sum_j p(y | M_j)p(M_j)}.$$

For **model selection**, consider posterior model odds:

$$\boxed{\frac{p(M_1 | y)}{p(M_2 | y)}} = \boxed{\frac{p(y | M_1)}{p(y | M_2)}} \times \boxed{\frac{p(M_1)}{p(M_2)}}.$$

posterior odds Bayes factor prior odds

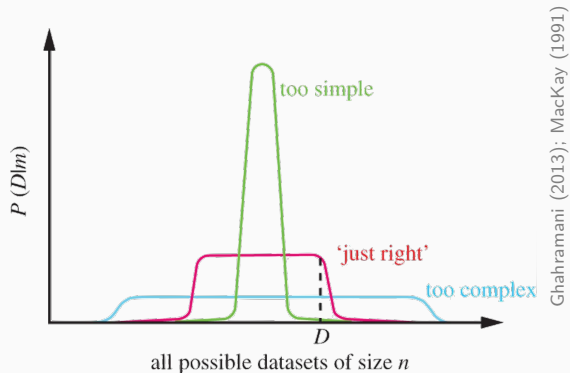
Must compute the **Bayesian model evidence** or **marginal likelihood** given by the normalising constant

$$z = p(y | M) = \int d\theta \mathcal{L}(\theta) \pi(\theta).$$

→ **Extremely challenging computational problem in high-dimensions.**

Occam's razor

The Bayesian model evidence **naturally incorporates Occam's razor**, trading off model complexity and goodness of fit.



On priors

- **Physics-informed priors**
e.g. mass constrained to be positive
- **Uninformative prior**
e.g. invariance to symmetry transformations
- **Informative prior**
e.g. regularize by imposing sparsity in dictionary
- **Data-informed priors**
e.g. prior \sim old data, likelihood \sim new data, posterior \sim old and new data
- **Data-driven priors**
e.g. empirical Bayes (estimate prior from data), learn by machine learning (generative models)

Challenge of Bayesian model selection

Naive Monte Carlo integration can be used to compute the marginal likelihood in principle.

However, the resulting estimator has very large variance, rendering it **ineffective in practice** (even in relatively low dimensional settings).

Challenge of Bayesian model selection

Naive Monte Carlo integration can be used to compute the marginal likelihood in principle.

However, the resulting estimator has very large variance, rendering it **ineffective in practice** (even in relatively low dimensional settings).

Require techniques **tailored** to the computation of the marginal likelihood.

Challenge of Bayesian model selection

Naive Monte Carlo integration can be used to compute the marginal likelihood in principle.

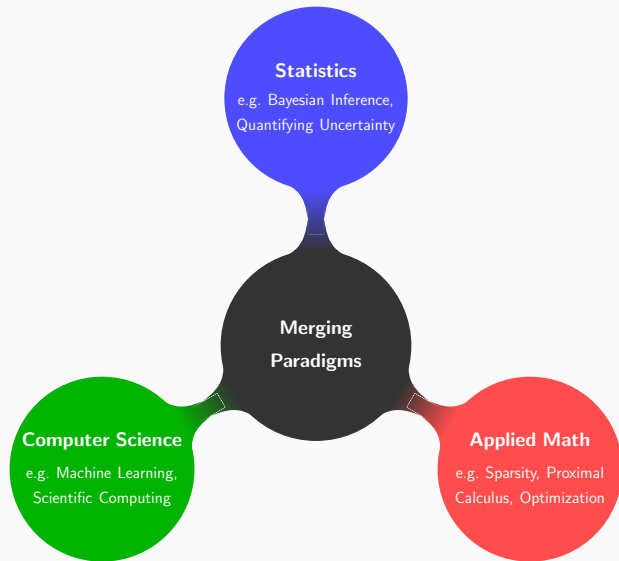
However, the resulting estimator has very large variance, rendering it **ineffective in practice** (even in relatively low dimensional settings).

Require techniques **tailored** to the computation of the marginal likelihood.

Challenges:

- Extending to **general sampling** strategies.
- Extending to **simulation-based inference** (likelihood-free inference).
- Scaling to **high-dimensions**.

Merging paradigms



Learnt harmonic mean estimator for likelihood-based model selection

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\rho = \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right]$$

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\rho = \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta|y)$$

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta|y) \\ &= \int d\theta \frac{1}{\mathcal{L}(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z}\end{aligned}$$

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta|y) \\ &= \int d\theta \frac{1}{\mathcal{L}(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} \\ &= \frac{1}{z}\end{aligned}$$

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta|y) \\ &= \int d\theta \frac{1}{\mathcal{L}(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} \\ &= \frac{1}{z}\end{aligned}$$

Original harmonic mean estimator (Newton & Raftery 1994)

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\mathcal{L}(\theta_i)}, \quad \theta_i \sim p(\theta|y)$$

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta|y) \\ &= \int d\theta \frac{1}{\mathcal{L}(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} \\ &= \frac{1}{z}\end{aligned}$$

Original harmonic mean estimator (Newton & Raftery 1994)

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\mathcal{L}(\theta_i)}, \quad \theta_i \sim p(\theta|y)$$

Very simple approach but **can fail catastrophically** (Neal 1994).

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\pi(\theta)}{p(\theta | y)} p(\theta | y) .$$

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\pi(\theta)}{p(\theta | y)} p(\theta | y) .$$

importance sampling

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\pi(\theta)}{p(\theta | y)} p(\theta | y) .$$

importance sampling

Importance sampling interpretation:

- Importance **sampling target distribution is prior** $\pi(\theta)$.
- Importance **sampling density is posterior** $p(\theta | y)$.

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\pi(\theta)}{p(\theta | y)} p(\theta | y) .$$

importance sampling

Importance sampling interpretation:

- Importance **sampling target distribution is prior** $\pi(\theta)$.
- Importance **sampling density is posterior** $p(\theta | y)$.

For importance sampling, want sampling density to have fatter tails than target.

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\pi(\theta)}{p(\theta | y)} p(\theta | y) .$$

importance sampling

Importance sampling interpretation:

- Importance **sampling target distribution is prior** $\pi(\theta)$.
- Importance **sampling density is posterior** $p(\theta | y)$.

For importance sampling, want sampling density to have fatter tails than target.

Not the case when importance sampling density is posterior and target is the prior.

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

***Re-targeted* harmonic mean relationship** (Gelfand & Dey 1994)

$$\rho = \mathbb{E}_{p(\theta|y)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right]$$

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

***Re-targeted* harmonic mean relationship** (Gelfand & Dey 1994)

$$\rho = \mathbb{E}_{p(\theta|y)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta|y)$$

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

Re-targeted harmonic mean relationship (Gelfand & Dey 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta|y) \\ &= \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z}\end{aligned}$$

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

Re-targeted harmonic mean relationship (Gelfand & Dey 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta|y) \\ &= \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} \\ &= \frac{1}{z}\end{aligned}$$

Re-targeted harmonic mean estimator

Introduce an arbitrary importance sampling target $\varphi(\theta)$ (which must be normalised).

Re-targeted harmonic mean relationship (Gelfand & Dey 1994)

$$\begin{aligned}\rho &= \mathbb{E}_{p(\theta|y)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta|y) \\ &= \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} \\ &= \frac{1}{z}\end{aligned}$$

Re-targeted harmonic mean estimator (Gelfand & Dey 1994)

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{\varphi(\theta_i)}{\mathcal{L}(\theta_i)\pi(\theta_i)}, \quad \theta_i \sim p(\theta|y)$$

Re-targeted harmonic mean estimator

Importance sampling interpretation:

$$\rho = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta|y) = \frac{1}{z} \int d\theta \frac{\varphi(\theta)}{p(\theta|y)} p(\theta|y).$$

Re-targeted harmonic mean estimator

Importance sampling interpretation:

$$\rho = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\varphi(\theta)}{p(\theta | y)} p(\theta | y).$$

Ensure importance sampling target $\varphi(\theta)$ does **not** have fatter tails than posterior $p(\theta | y)$ (importance sampling density).

Re-targeted harmonic mean estimator

Importance sampling interpretation:

$$\rho = \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} p(\theta | y) = \frac{1}{z} \int d\theta \frac{\varphi(\theta)}{p(\theta | y)} p(\theta | y).$$

Ensure importance sampling target $\varphi(\theta)$ does **not** have fatter tails than posterior $p(\theta | y)$ (importance sampling density).

→ **How set importance sampling target distribution $\varphi(\theta)$?**

How set importance sampling target distribution $\varphi(\theta)$?

Variety of cases been considered:

- Multi-variate Gaussian (e.g. Chib 1995)
- Indicator functions (e.g. Robert & Wraith 2009, van Haasteren 2009)

How set importance sampling target distribution $\varphi(\theta)$?

Variety of cases been considered:

- Multi-variate Gaussian (e.g. Chib 1995)
- Indicator functions (e.g. Robert & Wraith 2009, van Haasteren 2009)

Optimal target:

$$\varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}$$

(resulting estimator has zero variance).

How set importance sampling target distribution $\varphi(\theta)$?

Variety of cases been considered:

- Multi-variate Gaussian (e.g. Chib 1995)
- Indicator functions (e.g. Robert & Wraith 2009, van Haasteren 2009)

Optimal target:

$$\varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}$$

(resulting estimator has zero variance).

But clearly **not feasible** since requires knowledge of the evidence z (recall the target must be normalised) → **requires problem to have been solved already!**

Learnt harmonic mean estimator

Propose the *learnt* harmonic mean estimator
(McEwen, Wallis, Price, Docherty 2021; [arXiv:2111.12720](https://arxiv.org/abs/2111.12720)).



Chris Wallis



Matt Price



Matthew Docherty

Learnt harmonic mean estimator

Learn an approximation of the optimal target distribution:

$$\varphi(\theta) \stackrel{\text{ML}}{\simeq} \varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}$$

Learnt harmonic mean estimator

Learn an approximation of the optimal target distribution:

$$\varphi(\theta) \stackrel{\text{ML}}{\simeq} \varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}$$

- Approximation not required to be highly accurate.
- Must not have fatter tails than posterior.

Learnt harmonic mean estimator

Learn an approximation of the optimal target distribution:

$$\varphi(\theta) \stackrel{\text{ML}}{\simeq} \varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}$$

- Approximation not required to be highly accurate.
- Must not have fatter tails than posterior.

Also develop strategy to estimate the variance of the estimator, its variance, and other sanity checks.

Learning the target distribution

Consider a **variety of machine learning approaches**:

- Uniform hyper-ellipsoid
- Kernel Density Estimation (KDE)
- Modified Gaussian mixture model (MGMM)

Fit model by **minimising variance of resulting estimator**, while ensuring unbiased, with possible regularisation:

$$\min \hat{\sigma}^2 + \lambda R \quad \text{subject to} \quad \hat{\rho} = \hat{\mu}_1$$

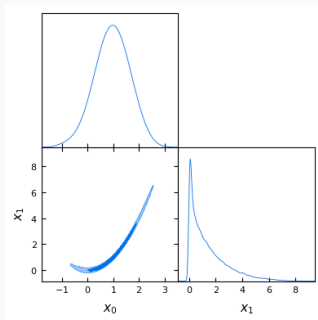
Solve by bespoke **mini-batch stochastic gradient descent**.

Cross-validation to select machine learning model and hyperparameters.

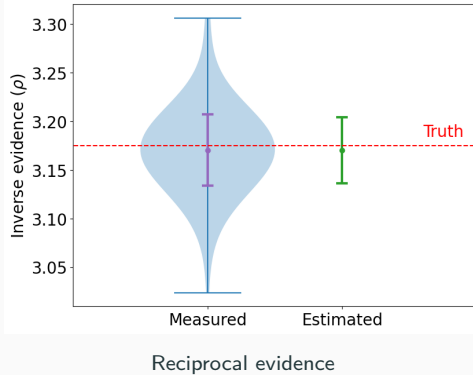
Rosenbrock example

Rosenbrock function is the classical example of a **pronounced thin curving degeneracy**, with likelihood defined by

$$f(\theta) = \sum_{i=1}^{n-1} \left[(a - \theta_i)^2 + b(\theta_{i+1} - \theta_i^2)^2 \right], \quad \log(\mathcal{L}(\theta)) = -f(\theta).$$

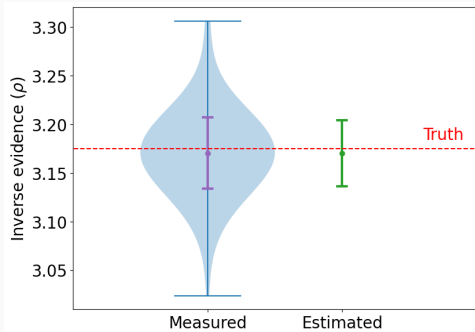


Rosenbrock example

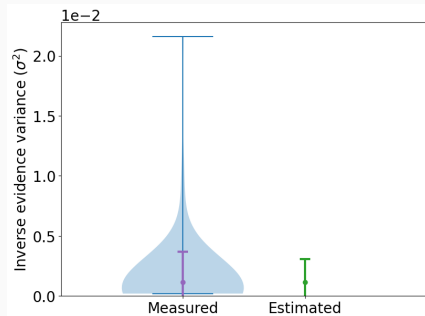


Accuracy of learnt harmonic mean estimator for Rosenbrock example.

Rosenbrock example



Reciprocal evidence



Variance of reciprocal evidence

Accuracy of learnt harmonic mean estimator for Rosenbrock example.

Normal-Gamma example

Pathological example (Friel & Wyse 2012) where original harmonic mean estimator fails.

Normal-Gamma example

Pathological example (Friel & Wyse 2012) where original harmonic mean estimator fails.

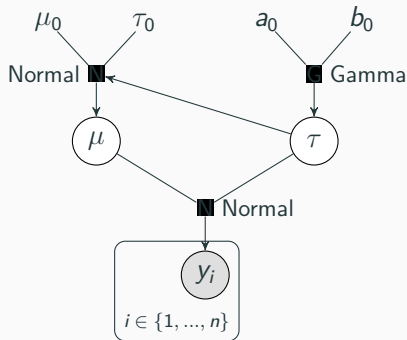
Data model:

$$y_i \sim N(\mu, \tau^{-1})$$

Prior model:

$$\text{Mean: } \mu \sim N(\mu_0, (\tau_0 \tau)^{-1})$$

$$\text{Precision: } \tau \sim \text{Ga}(a_0, b_0)$$



Normal-Gamma example

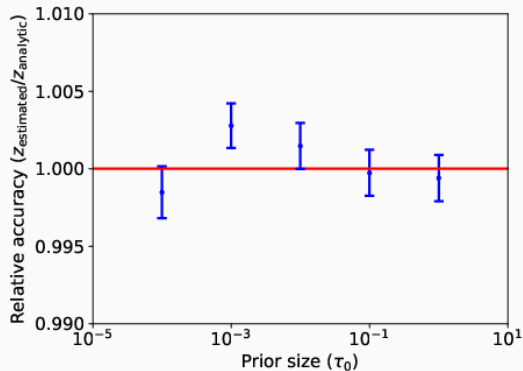
Analytic evidence:

$$z = (2\pi)^{-n/2} \frac{\Gamma(a_n) b_0^{a_0}}{\Gamma(a_0) b_n^{a_n}} \left(\frac{\tau_0}{\tau_n} \right)^{1/2}$$

where

$$\tau_n = \tau_0 + n, \quad a_n = a_0 + n/2, \quad b_n = b_0 + \frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^2 + \frac{\tau_0 n (\bar{y} - \mu_0)^2}{2(\tau_0 + n)}.$$

Normal-Gamma example



Comparison of marginal likelihood values computed to truth for varying prior.

Normal-Gamma example

Marginal likelihood values for Normal-Gamma example with varying prior.

τ_0	10^{-4}	10^{-3}	10^{-2}	10^{-1}	10^0
Analytic $\log(z)$	-144.5530	-143.4017	-142.2505	-141.0999	-139.9552
Estimated $\log(\hat{z})$	-144.5545	-143.3990	-142.2490	-141.1001	-139.9558
Error (learnt harmonic mean)	-0.0015	0.0027	0.0015	-0.0011	-0.0006
Error (original harmonic mean)	12.2100	—	9.7900	8.5000	7.1000

Radiata pine example

Radiata pine data-set has become **classical benchmark** for evaluating evidence estimators:

- maximum compression strength parallel to grain y_i ,
- density x_i ,
- density adjust for resin content z_i ,

for $i \in \{1, \dots, n\}$ where $n = 42$ specimens.



Is **density** or **resin-adjusted density** a better predictor of compression strength?

Radiata pine example

Gaussian linear models:

$$M_1 : \quad y_i = \alpha + \underbrace{\beta(x_i - \bar{x})}_{\text{density}} + \epsilon_i, \quad \epsilon_i \sim \text{N}(0, \tau^{-1}).$$

$$M_2 : \quad y_i = \gamma + \underbrace{\delta(z_i - \bar{z})}_{\text{resin-adjusted density}} + \eta_i, \quad \eta_i \sim \text{N}(0, \lambda^{-1}).$$

Priors for model 1 (similar for model 2):

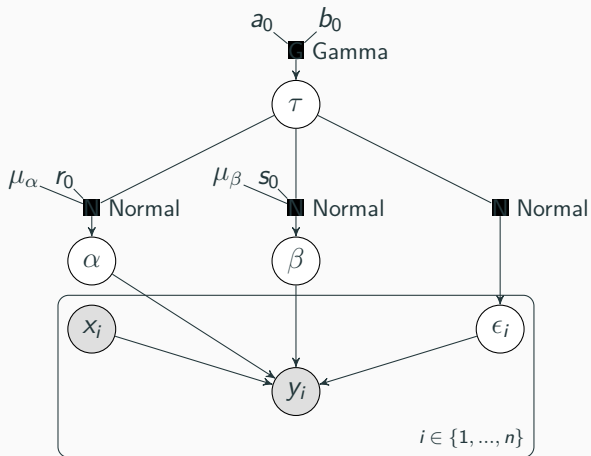
$$\alpha \sim \text{N}(\mu_\alpha, (r_0\tau)^{-1}),$$

$$\beta \sim \text{N}(\mu_\beta, (s_0\tau)^{-1}),$$

$$\tau \sim \text{Ga}(a_0, b_0),$$

$$(\mu_\alpha = 3000, \mu_\beta = 185, r_0 = 0.06, s_0 = 6, a_0 = 3, b_0 = 2 \times 300^2).$$

Radiata pine example



Hierarchical Bayesian model for Radiata pine example (for model 1; model 2 is similar).

Radiata pine example

Analytic evidence:

$$z = \pi^{-n/2} b_0^{a_0} \frac{\Gamma(a_0 + n/2) |Q_0|^{1/2}}{\Gamma(a_0) |M|^{1/2}} (\mathbf{y}^T \mathbf{y} + \boldsymbol{\mu}_0^T Q_0 \boldsymbol{\mu}_0 - \boldsymbol{\nu}_0^T M \boldsymbol{\nu}_0 + 2b_0)^{-a_0 - n/2}$$

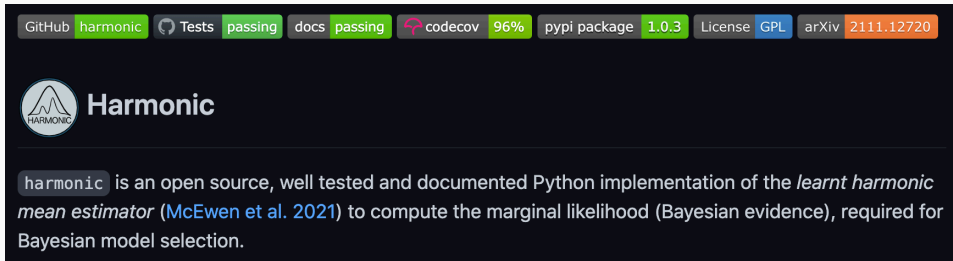
where $\boldsymbol{\mu}_0 = (\mu_\alpha, \mu_\beta)^T$, $Q_0 = \text{diag}(r_0, s_0)$, and $M = X^T X + Q_0$.

Radiata pine example


Marginal likelihood values for Radiata Pine example.

	Model M_1 $\log(z_1)$	Model M_2 $\log(z_2)$	$\log \text{BF}_{21}$ $= \log(z_2) - \log(z_1)$
Analytic	-310.12829	-301.70460	8.42368
Estimated	-310.12807	-301.70413	8.42394
	± 0.00072	± 0.00074	± 0.00145
Error (learnt harmonic mean)	0.00022	0.00047	0.00026
Error (original harmonic mean)	–	–	-0.17372

Harmonic code



GitHub `harmonic` Tests `passing` docs `passing` codecov `96%` pypi package `1.0.3` License `GPL` arXiv `2111.12720`

 **Harmonic**

`harmonic` is an open source, well tested and documented Python implementation of the *learnt harmonic mean estimator* (McEwen et al. 2021) to compute the marginal likelihood (Bayesian evidence), required for Bayesian model selection.

Github: `https://github.com/astro-informatics/harmonic`

Docs: `https://astro-informatics.github.io/harmonic`

(Seamless integration with emcee.)

Learnt harmonic mean estimator for simulation-based model selection

Simulation-based inference (SBI)

Consider situation where the likelihood $p(y|\theta, M)$ is unknown or intractable.

Simulation-based inference (likelihood-free inference) seeks to perform parameter inference by estimating the posterior $p(\theta | y_o, M)$ for observed data y_o using **simulations only**.

Simulation-based inference (SBI)

Consider situation where the likelihood $p(y|\theta, M)$ is unknown or intractable.

Simulation-based inference (likelihood-free inference) seeks to perform parameter inference by estimating the posterior $p(\theta | y_o, M)$ for observed data y_o using **simulations only**.

Advantages:

- Forward modelling of complex physics, contamination, observational process.
- No assumptions on the form of the likelihood.

Sequential neural posterior estimation (SNPE)

SNPE introduced by Papamakarios & Murray (2016).

Construct training data $\{(\theta_i, y_i)\}$ where parameter drawn from proposal prior $\theta_i \sim \tilde{p}(\theta | M)$ and then generate simulation $y_i \sim p(y | \theta_i)$.

Learn posterior

$$q_\phi(\theta | y, M) \simeq p(\theta | y, M),$$

where ϕ are the parameters of the learned model.

Sequential neural posterior estimation (SNPE)

SNPE introduced by Papamakarios & Murray (2016).

Construct training data $\{(\theta_i, y_i)\}$ where parameter drawn from proposal prior $\theta_i \sim \tilde{p}(\theta | M)$ and then generate simulation $y_i \sim p(y | \theta_i)$.

Learn posterior

$$q_\phi(\theta | y, M) \simeq p(\theta | y, M),$$

where ϕ are the parameters of the learned model.

Train by maximising the probability of the parameters vectors $\prod_i q_\phi(\theta_i | y_i, M)$, *i.e.* minimising loss function

$$L(\phi) = - \sum_i \log q_\phi(\theta_i | y_i, M).$$

Efficient proposal prior

If proposal prior matches underlying prior, *i.e.* $\tilde{p}(\theta | M) = p(\theta | M)$, then would learn an unbiased estimate of the posterior.

BUT, this **highly inefficient** since learn posterior for all possible data realisations but only interested in observational data y_o .

Efficient proposal prior

If proposal prior matches underlying prior, *i.e.* $\tilde{p}(\theta | M) = p(\theta | M)$, then would learn an unbiased estimate of the posterior.

BUT, this **highly inefficient** since learn posterior for all possible data realisations but only interested in observational data y_o .

Sequential approach: Train in *runs* where the proposal prior matches the intermediate learned posterior at y_o . Update proposal prior for each run.

Efficient proposal prior

If proposal prior matches underlying prior, *i.e.* $\tilde{p}(\theta | M) = p(\theta | M)$, then would learn an unbiased estimate of the posterior.

BUT, this **highly inefficient** since learn posterior for all possible data realisations but only interested in observational data y_o .

Sequential approach: Train in *runs* where the proposal prior matches the intermediate learned posterior at y_o . Update proposal prior for each run.

Since do not sample from prior, end up learning

$$q_{\phi}(\theta | y) \simeq \frac{\tilde{p}(\theta | M)}{p(\theta | M)} p(\theta | y = y_o, M) .$$

Correcting the bias of SNPE

Variants of SNPE introduced to **correct the bias** of sampling from proposal prior.

Correcting the bias of SNPE

Variants of SNPE introduced to **correct the bias** of sampling from proposal prior.

- ▷ SNPE-A (Papamakarios & Murray 2016): importance weight learned distribution.
- ▷ SNPE-B (Lueckmann *et al.* 2017): adjust loss function to include importance weights.
- ▷ SNPE-C (Greenberg *et al.* 2019): reparameterise the proposal posterior objective to recover the true posterior distribution.

Sequential neural likelihood estimation (SNLE)

Avoid the bias issue of SNPE by learning the likelihood \Rightarrow **SNLE** (Papamakarios *et al.* 2019).

Sequential neural likelihood estimation (SNLE)

Avoid the bias issue of SNPE by learning the likelihood \Rightarrow **SNLE** (Papamakarios *et al.* 2019).

Define joint distribution $\tilde{p}(\theta, y) = p(y|\theta, M)\tilde{p}(\theta, M)$.

Learn likelihood

$$q_{\psi}(y|\theta, M) \simeq p(y|\theta, M),$$

where ψ are the parameters of the learned model.

Sequential neural likelihood estimation (SNLE)

Avoid the bias issue of SNPE by learning the likelihood \Rightarrow **SNLE** (Papamakarios *et al.* 2019).

Define joint distribution $\tilde{p}(\theta, y) = p(y|\theta, M)\tilde{p}(\theta, M)$.

Learn likelihood

$$q_{\psi}(y|\theta, M) \simeq p(y|\theta, M),$$

where ψ are the parameters of the learned model.

Train by maximising total log-likelihood $\sum_i \log q_{\psi}(y_i|\theta_i)$, *i.e.* by maximising

$$\mathbb{E}_{\tilde{p}(\theta, y)}[\log q_{\psi}(y|\theta, M)] = -\mathbb{E}_{\tilde{p}(\theta)}[D_{\text{KL}}(p(y|\theta, M), q_{\psi}(y|\theta, M))] + \text{const.},$$

where D_{KL} is the Kullback-Leibler divergence.

Sequential approach and sampling

Adopt sequential approach for SNLE, where consider run r and update proposal prior $\tilde{p}(\theta | M)$ to use current posterior estimate, *i.e.* $p_r(\theta | y_o, M) \propto q_{\psi}(y_o | \theta, M)p(\theta | M)$.

Unlike, SNPE, **no adjustment is necessary to account for proposing strategy.**

Sample from approximate posterior by MCMC sampling.

Bayesian model comparison for simulation-based inference

Bayesian model comparison for simulation-based inference

(Spurio Mancini, Docherty, Price, McEwen, in prep.).



Alessio Spurio Mancini



Matthew Docherty



Matt Price

Naive model evidence computation

Recall SNPE and SNLE:

$$q_{\phi}^{\text{SNPE}}(\theta | y, M) \simeq p(\theta | y, M); \quad q_{\psi}^{\text{SNLE}}(y | \theta, M) \simeq p(y | \theta, M).$$

Naive estimate of the model evidence:

$$\hat{z} = \frac{1}{N} \sum_i \frac{q_{\psi}^{\text{SNLE}}(y_o | \theta_i, M) p(\theta_i | M)}{q_{\phi}^{\text{SNPE}}(\theta_i | y_o, M)}.$$

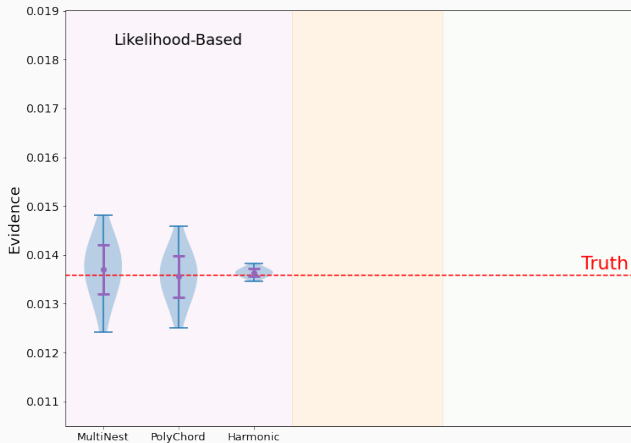
Computing model evidence by learnt harmonic mean estimator

Compute model evidence by learnt harmonic mean estimator:

$$\hat{\rho} = \frac{1}{N} \sum_i \frac{\varphi(\theta_i)}{q_{\psi}^{\text{SNLE}}(y_o | \theta_i, M) p(\theta_i | M)}.$$

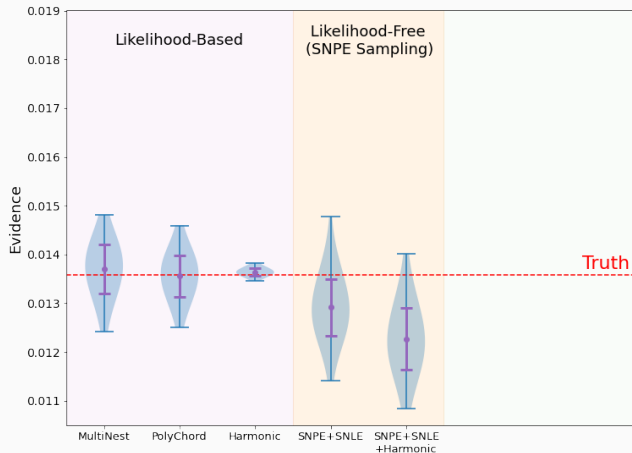
- Agnostic to sampling strategy.
- If adopt SNLE only, then sample by MCMC. If adopt SNPE for samples, then still require SNLE for likelihood.

Linear Gaussian example



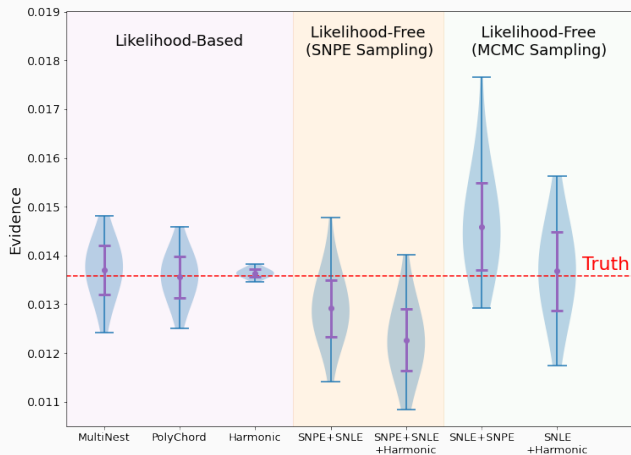
Model evidence computed in likelihood-based and likelihood-free settings.

Linear Gaussian example



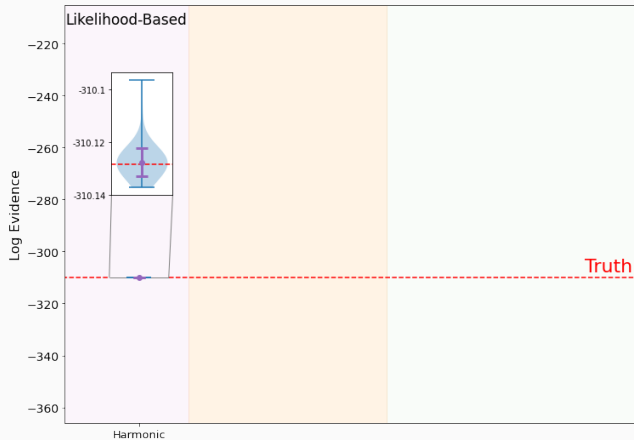
Model evidence computed in likelihood-based and likelihood-free settings.

Linear Gaussian example



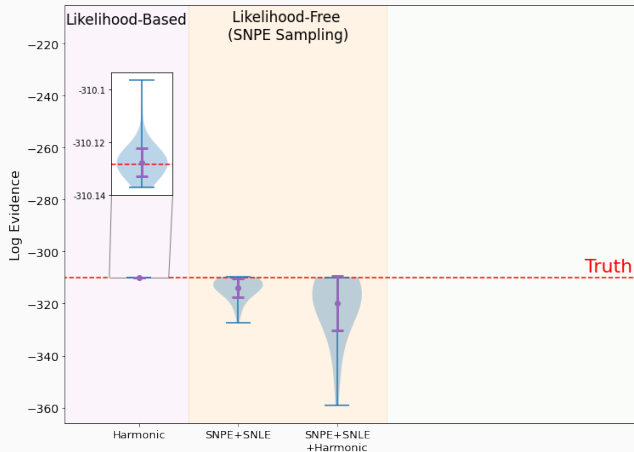
Model evidence computed in likelihood-based and likelihood-free settings.

Radiata pine example



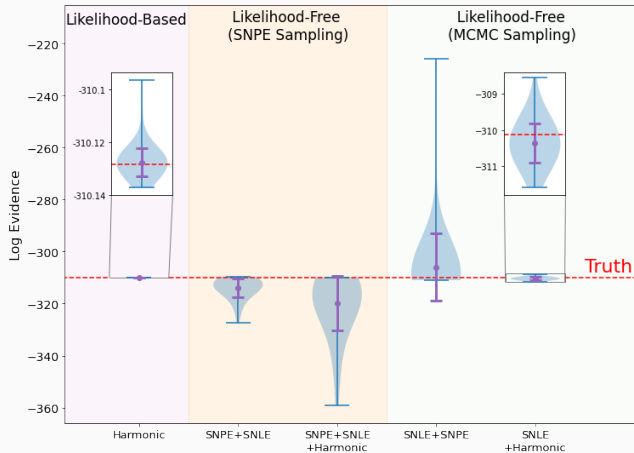
Model evidence computed in likelihood-based and likelihood-free settings.

Radiata pine example



Model evidence computed in likelihood-based and likelihood-free settings.

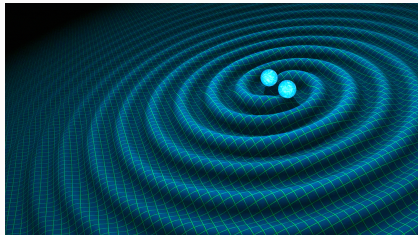
Radiata pine example



Model evidence computed in likelihood-based and likelihood-free settings.

Gravitational waves

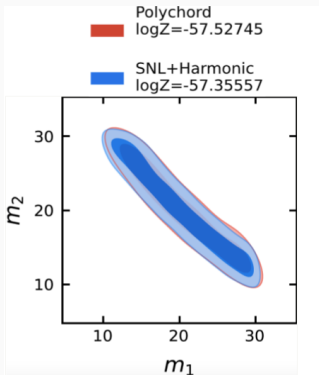
Simulate a **black-hole, black-hole merger** as observed by an interferometer (e.g. LIGO).



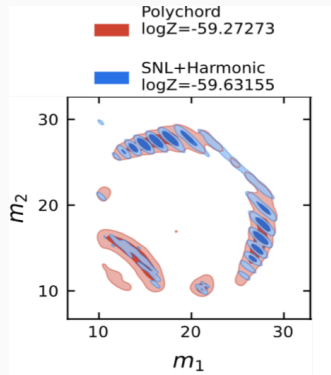
Consider two models (*Spin-Precessing Effective-One-Body Numerical Relativity* and *Inspiral Ringdown Merger*) and perform model comparison.

Likelihood available for validation.

Gravitational waves



*Spin-Precessing
Effective-One-Body Numerical
Relativity
(correct model)*



*Inspiral Ringdown Merger
(incorrect model)*

Proximal nested sampling for high-dimensional model selection

Nested sampling: reparameterising the likelihood

Nested sampling is a clever approach to efficiently evaluate the evidence (Skilling 2006).

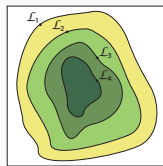
Consider $\Omega_{L^*} = \{x | \mathcal{L}(x) \geq L^*\}$, which groups the parameter space Ω into a series of **nested subspaces**.

Define the prior volume ξ within Ω_{L^*} by $\xi(L^*) = \int_{\Omega_{L^*}} \pi(x) dx$.

The marginal likelihood integral can then be rewritten as

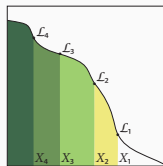
$$\mathcal{Z} = \int_0^1 \mathcal{L}(\xi) d\xi,$$

which is a **one-dimensional integral** over the prior volume ξ .



Feroz et al. (2013)

Nested subspaces



Feroz et al. (2013)

Reparameterised
likelihood

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

1. Draw N_{live} *live* samples from prior, with prior volume $\xi_0 = 1$.

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

1. Draw N_{live} *live* samples from prior, with prior volume $\xi_0 = 1$.
2. Remove sample with smallest likelihood, say L_i .

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

1. Draw N_{live} *live* samples from prior, with prior volume $\xi_0 = 1$.
2. Remove sample with smallest likelihood, say L_i .
3. Replace removed sample with new sample from the prior but *constrained to a higher likelihood* than L_i .

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

1. Draw N_{live} *live* samples from prior, with prior volume $\xi_0 = 1$.
2. Remove sample with smallest likelihood, say L_i .
3. Replace removed sample with new sample from the prior but *constrained to a higher likelihood* than L_i .
4. Estimate (stochastically) prior volume ξ_i enclosed by likelihood level-set L_i .

Nested sampling: constrained sampling

Require strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Nested sampling (Skilling 2006)

1. Draw N_{live} *live* samples from prior, with prior volume $\xi_0 = 1$.
2. Remove sample with smallest likelihood, say L_i .
3. Replace removed sample with new sample from the prior but *constrained to a higher likelihood* than L_i .
4. Estimate (stochastically) prior volume ξ_i enclosed by likelihood level-set L_i .
5. Repeat 2–5.

Nested sampling: estimating enclosed prior volume stochastically

Enclosed prior volume decreases exponentially at each step: $\xi_{i+1} = t_{i+1}\xi_i$.

Shrinkage ratio can be estimated stochastically since $\mathbb{E}(\log t) = -1/N_{\text{live}}$.

The enclosed prior volume can then be estimated by

$$\xi_{i+1} = \exp(-i/N_{\text{live}}).$$

Nested sampling: evidence estimation and posterior inference

Given the sequence of decreasing prior volumes $\{\xi_i\}_{i=0}^N$ and corresponding likelihoods $L_i = \mathcal{L}(\xi_i)$, the **model evidence** can be computed numerically using standard quadrature:

$$\mathcal{Z} = \sum_{i=1}^N L_i w_i ,$$

for quadrature weight w_i (e.g. the trapezium rule with $w_i = (\xi_{i-1} + \xi_{i+1})/2$).

Nested sampling: evidence estimation and posterior inference

Given the sequence of decreasing prior volumes $\{\xi_i\}_{i=0}^N$ and corresponding likelihoods $L_i = \mathcal{L}(\xi_i)$, the **model evidence** can be computed numerically using standard quadrature:

$$\mathcal{Z} = \sum_{i=1}^N L_i w_i ,$$

for quadrature weight w_i (e.g. the trapezium rule with $w_i = (\xi_{i-1} + \xi_{i+1})/2$).

Posterior inferences can also be computed by assigning importance weights

$$p_i = \frac{L_i w_i}{\mathcal{Z}} .$$

Nested sampling: challenge

Recall: to compute the marginal likelihood by nested sampling require strategy to generate likelihoods L_i and associated prior volumes ξ_i .

Achieved by **sampling from the prior, subject the likelihood iso-contour constraint**, *i.e.* sampling from the prior $\pi(x)$, such that $\mathcal{L}(x) > L^*$

Nested sampling: challenge

Recall: to compute the marginal likelihood by nested sampling require strategy to generate likelihoods L_i and associated prior volumes ξ_i .

Achieved by **sampling from the prior, subject the likelihood iso-contour constraint**, *i.e.* sampling from the prior $\pi(x)$, such that $\mathcal{L}(x) > L^*$

This is the **main difficulty** in applying nested sampling to high-dimensional problems.

Exploit common structure

Many high-dimensional inverse problems are **log-convex**, e.g. inverse imaging problems with Gaussian data fidelity and sparsity-promoting prior.

Exploit structure (log convexity) of the problem.

⇒ **Proximal nested sampling** (Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](https://arxiv.org/abs/2106.03646))



Xiaohao Cai



Marcelo Pereyra

Constrained sampling formulation

Consider case where prior and likelihood of form

$$\pi(x) = \exp(-f(x)),$$

prior

$$\mathcal{L}(x) = \exp(-g(x)),$$

likelihood

where f and g are convex lower semicontinuous functions on Ω .

Let $\iota_{L^*}(x)$ and $\chi_{L^*}(x)$ be the indicator and characteristic functions:

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

Then let $\pi_{L^*}(x) = \pi(x)\iota_{L^*}(x)$ represent the prior distribution with the hard likelihood constraint.

Constrained sampling formulation

Taking the logarithm, we can write

$$-\log \pi_{L^*}(x) = f(x) + \chi_{\mathcal{B}_\tau}(x),$$

where $\chi_{\mathcal{B}_\tau}(x)$ is the characteristic function associated with the convex set

$$\mathcal{B}_\tau := \{x | g(x) < \tau\},$$

for $\tau = -\log L^*$.

MCMC sampling with Langevin dynamics

Consider posteriors of the following form:

$$p(\mathbf{x} | \mathbf{y}) = \pi(\mathbf{x}) \propto \exp(-p(\mathbf{x})).$$

If $p(\mathbf{x})$ differentiable can adopt Langevin dynamics.

Based on **Langevin diffusion process** $\mathcal{L}(t)$, with π as stationary distribution:

$$d\mathcal{L}(t) = \frac{1}{2} \nabla \log \pi(\mathcal{L}(t)) dt + d\mathcal{W}(t), \quad \mathcal{L}(0) = l_0$$

where \mathcal{W} is Brownian motion.

MCMC sampling with Langevin dynamics

Consider posteriors of the following form:

$$p(\mathbf{x} | \mathbf{y}) = \pi(\mathbf{x}) \propto \exp(-p(\mathbf{x})).$$

If $p(\mathbf{x})$ differentiable can adopt Langevin dynamics.

Based on **Langevin diffusion process** $\mathcal{L}(t)$, with π as stationary distribution:

$$d\mathcal{L}(t) = \frac{1}{2} \underbrace{\nabla \log \pi(\mathcal{L}(t))}_{\text{gradient}} dt + d\mathcal{W}(t), \quad \mathcal{L}(0) = l_0$$

where \mathcal{W} is Brownian motion.

Need gradients so **not directly applicable**.

Moreau-Yosida approximation

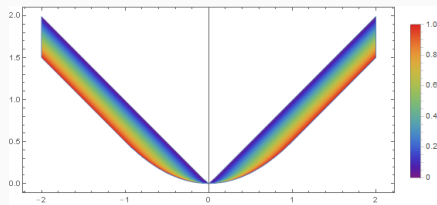
Moreau-Yosida approximation (envelope) of

f :

$$\hat{f}^\lambda(\mathbf{x}) = \inf_{\mathbf{u} \in \mathbb{R}^N} f(\mathbf{u}) + \frac{\|\mathbf{u} - \mathbf{x}\|^2}{2\lambda}$$

Important properties of $\hat{f}^\lambda(\mathbf{x})$:

1. As $\lambda \rightarrow 0$, $\hat{f}^\lambda(\mathbf{x}) \rightarrow f(\mathbf{x})$
2. $\nabla \hat{f}^\lambda(\mathbf{x}) = (\mathbf{x} - \text{prox}_f^\lambda(\mathbf{x}))/\lambda$



Moreau-Yosida envelope of $|x|$ for varying λ
[Credit: Stack exchange (ubpdqn)]

Proximal nested sampling

Proximal nested sampling (Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](https://arxiv.org/abs/2106.03646))

- Constrained sampling formulation
- Langevin MCMC sampling
- Moreau-Yosida approximation of constraint (and any non-differentiable prior)

Proximal nested sampling

Proximal nested sampling (Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](https://arxiv.org/abs/2106.03646))

- Constrained sampling formulation
- Langevin MCMC sampling
- Moreau-Yosida approximation of constraint (and any non-differentiable prior)

Proximal nested sampling Markov chain:

$$x^{(k+1)} = x^{(k)} - \frac{\delta}{2} \nabla f(x^{(k)}) - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{\chi_{B_\tau}}(x^{(k)})] + \sqrt{\delta} w^{(k+1)} .$$

Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

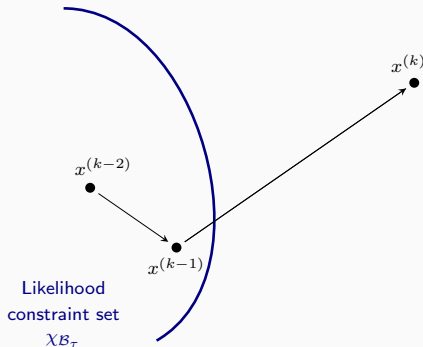
$$x^{(k+1)} = x^{(k)} - \frac{\delta}{2} \nabla f(x^{(k)}) - \frac{\delta}{2\lambda} [x^{(k)} - \text{prox}_{\chi_{B_\tau}}(x^{(k)})] + \sqrt{\delta} w^{(k+1)}.$$

Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

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

1. $x^{(k)}$ is already in \mathcal{B}_τ : term $[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$ disappears and recover usual Langevin MCMC.

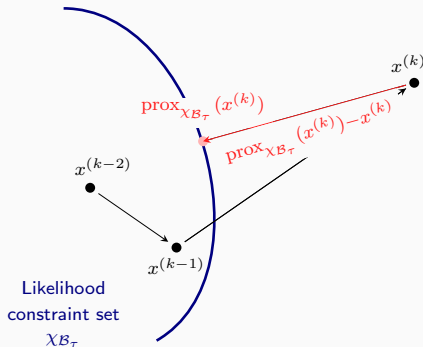


Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

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

1. $x^{(k)}$ **is already in** \mathcal{B}_τ : term $[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$ disappears and recover usual Langevin MCMC.
2. $x^{(k)}$ **is not in** \mathcal{B}_τ : a step is also taken in the direction $-[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$, which moves the next iteration in the direction of the projection of $x^{(k)}$ onto the convex set \mathcal{B}_τ . Acts to push the Markov chain back into the constraint set \mathcal{B}_τ if it wanders outside of it.

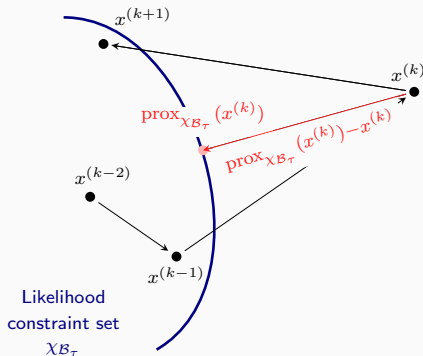


Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

$$x^{(k+1)} = x^{(k)} - \frac{\delta}{2} \nabla f(x^{(k)}) - \frac{\delta}{2\lambda} \left[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)}) \right] + \sqrt{\delta} w^{(k+1)}.$$

1. $x^{(k)}$ **is already in** \mathcal{B}_τ : term $[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$ disappears and recover usual Langevin MCMC.
2. $x^{(k)}$ **is not in** \mathcal{B}_τ : a step is also taken in the direction $-[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$, which moves the next iteration in the direction of the projection of $x^{(k)}$ onto the convex set \mathcal{B}_τ . Acts to push the Markov chain back into the constraint set \mathcal{B}_τ if it wanders outside of it.



Proximal nested sampling

A subsequent **Metropolis-Hastings** step guarantees hard likelihood constraint is satisfied.

Proximal nested sampling

A subsequent **Metropolis-Hastings** step guarantees hard likelihood constraint is satisfied.

In practice need to compute $\text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})$, including measurement operator.

For sparsity-promoting non-differentiable priors $f(x)$, can also make Moreau-Yosida approximation $f^\lambda(x)$ and leverage prox to compute gradient ∇f^λ .

Proximal nested sampling

A subsequent **Metropolis-Hastings** step guarantees hard likelihood constraint is satisfied.

In practice need to compute $\text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})$, including measurement operator.

For sparsity-promoting non-differentiable priors $f(x)$, can also make Moreau-Yosida approximation $f^\lambda(x)$ and leverage prox to compute gradient ∇f^λ .

Many further details regarding **explicit forms of proximal nested sampling** for common priors and likelihoods and how to compute proximity operators efficiently (Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](https://arxiv.org/abs/2106.03646)).

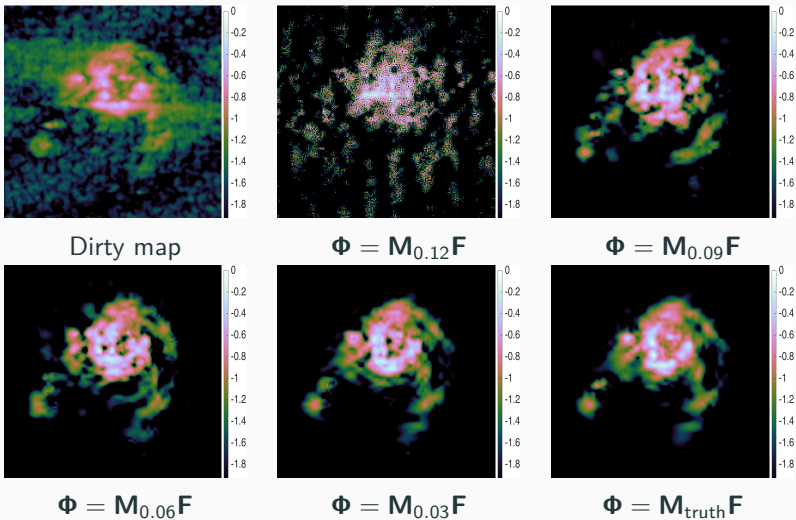
Measurement model misspecification experiment

Consider ground truth model $\Phi = \mathbf{M}_{\text{truth}}\mathbf{F}$ to simulate observational data \mathbf{y} .

However, when solving the inverse problem consider misspecified models \mathbf{M}_γ , where $\gamma > 0$ encodes the level of misspecification (mimics incorrectly specified wavelength).

Compute the model evidence using **proximal nested sampling**, using evidence to distinguish correct model.

Measurement model misspecification experiment



Measurement model misspecification experiment

Model	$\log \mathcal{Z}$	RMSE (Requires ground truth)
$\Phi = \mathbf{M}_{\text{truth}}\mathbf{F}$	$-4.47 \times 10^3 \pm 0.08$	3.40
$\Phi = \mathbf{M}_{0.03}\mathbf{F}$	$-4.88 \times 10^3 \pm 0.08$	7.85
$\Phi = \mathbf{M}_{0.06}\mathbf{F}$	$-5.63 \times 10^3 \pm 0.08$	12.01
$\Phi = \mathbf{M}_{0.09}\mathbf{F}$	$-9.21 \times 10^3 \pm 0.07$	15.71
$\Phi = \mathbf{M}_{0.12}\mathbf{F}$	$-1.44 \times 10^4 \pm 0.08$	18.08

Evidence computed by proximal nested sampling correctly classifies models.

Summary

Summary

Many science questions are **questions of model comparison**.

In cosmology we have only one Universe to observe \Rightarrow **Bayesian model selection**.

Many outstanding challenges:

- Extending to **general sampling** strategies.
 - Extending to **simulation-based inference** (likelihood-free inference).
 - Scaling to **high-dimensions**.
 - Learned **data-driven priors**.
1. Learnt harmonic mean estimator for Bayesian model comparison (McEwen, Wallis, Price, Docherty 2021; [arXiv:2111.12720](#))
 2. Bayesian model comparison for simulation-based inference (Spurio Mancini, Docherty, Price, McEwen, in prep.).
 3. Proximal nested sampling for high-dimensional Bayesian model comparison (Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](#))