

# Bayesian uncertainty quantification

for radio interferometry and beyond

---

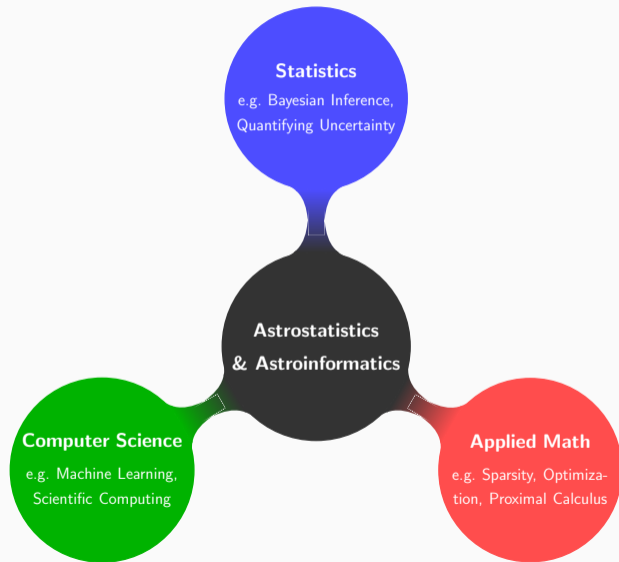
Jason D. McEwen

[www.jasonmcewen.org](http://www.jasonmcewen.org)

Mullard Space Science Laboratory (MSSL), UCL

April 2022

# Merging paradigms



# Bayesian inference: parameter estimation

## Bayes' theorem

$$\underbrace{P(\theta | y, M)}_{\text{posterior}} = \frac{\underbrace{P(y | \theta, M)}_{\text{likelihood}} \underbrace{P(\theta | M)}_{\text{prior}}}{\underbrace{P(y | M)}_{\text{evidence}}},$$

for parameters  $\theta$ , model  $M$  and observed data  $y$ .

# Bayesian inference: parameter estimation

## Bayes' theorem

$$\underbrace{P(\theta | y, M)}_{\text{posterior}} = \frac{\underbrace{P(y | \theta, M)}_{\text{likelihood}} \underbrace{P(\theta | M)}_{\text{prior}}}{\underbrace{P(y | M)}_{\text{evidence}}},$$

for parameters  $\theta$ , model  $M$  and observed data  $y$ .

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

# Bayesian inference: parameter estimation

## Bayes' theorem

$$\underbrace{P(\theta | y, M)}_{\text{posterior}} = \frac{\underbrace{P(y | \theta, M)}_{\text{likelihood}} \underbrace{P(\theta | M)}_{\text{prior}}}{\underbrace{P(y | M)}_{\text{evidence}}},$$

for parameters  $\theta$ , model  $M$  and observed data  $y$ .

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

→ **Challenging computational problem in high-dimensions.**

# Bayesian inference: model selection

For model selection, consider the posterior model probabilities:

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

posterior odds      prior odds      Bayes factor

# Bayesian inference: model selection

For **model selection**, consider the posterior model probabilities:

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

posterior odds      prior odds      Bayes factor

Must compute the **Bayesian 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

For **model selection**, consider the posterior model probabilities:

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

posterior odds      prior odds      Bayes factor

Must compute the **Bayesian 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.**



1. Learnt harmonic mean estimator for Bayesian model comparison
2. Proximal nested sampling for high-dimensional Bayesian model comparison
3. High-dimensional Bayesian uncertainty quantification for extreme computation

## Learnt harmonic mean estimator for Bayesian model comparison

---

# Desirable properties for Bayesian evidence estimators

Seek estimator that is:

- **Agnostic to sampling method** and uses posterior samples.
- Potential to **scale to high-dimensions**.

# Desirable properties for Bayesian evidence estimators

Seek estimator that is:

- **Agnostic to sampling method** and uses posterior samples.
- Potential to **scale to high-dimensions**.

*Harmonic mean estimator* has potential to meet these criteria but has serious shortcomings as originally posed.

# 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).

Recall:

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

# 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).

Recall:

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

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 *et al.* 2021; [arXiv:2111.12720](https://arxiv.org/abs/2111.12720)).

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

Propose the **learnt harmonic mean estimator** (McEwen *et al.* 2021; [arXiv:2111.12720](https://arxiv.org/abs/2111.12720)).

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

Propose the **learnt harmonic mean estimator** (McEwen *et al.* 2021; [arXiv:2111.12720](https://arxiv.org/abs/2111.12720)).

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**.

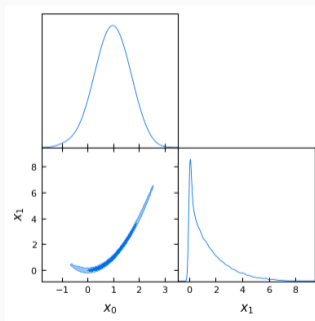
**Cross-validation** to select machine learning model and hyperparameters.



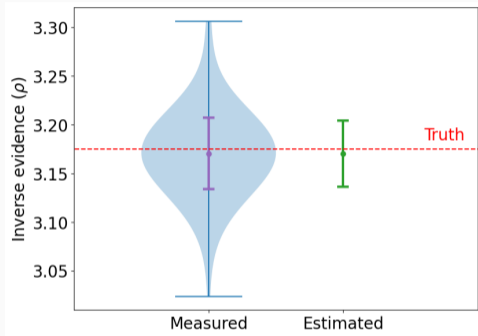
## Rosenbrock example

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

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



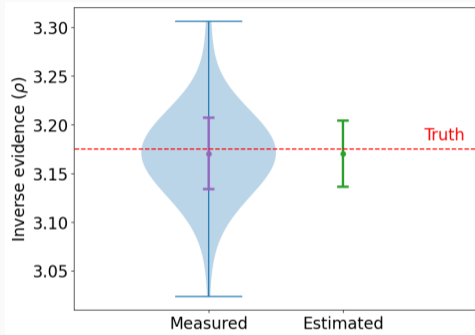
# Rosenbrock example



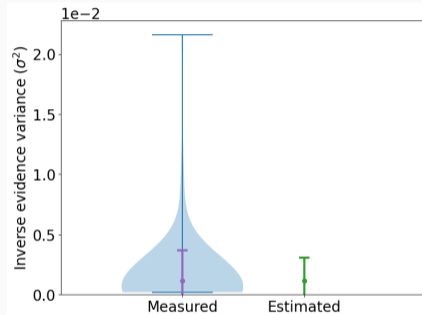
Reciprocal evidence

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.

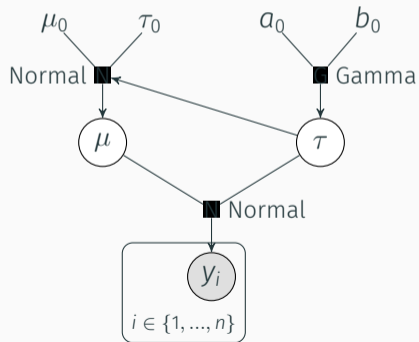
Data model:

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

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

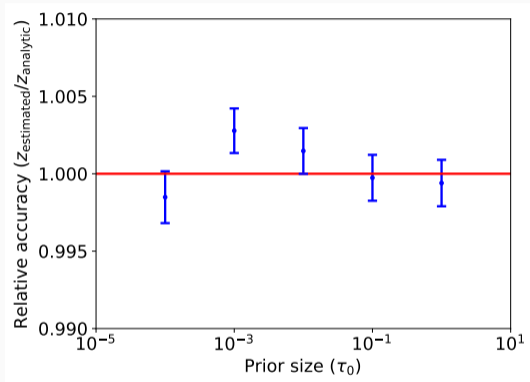
Analytic evidence:

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

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.

$\tau_0$	$10^{-4}$	$10^{-3}$	$10^{-2}$	$10^{-1}$	$10^0$
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:

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

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

Priors for model 1 (similar for model 2):

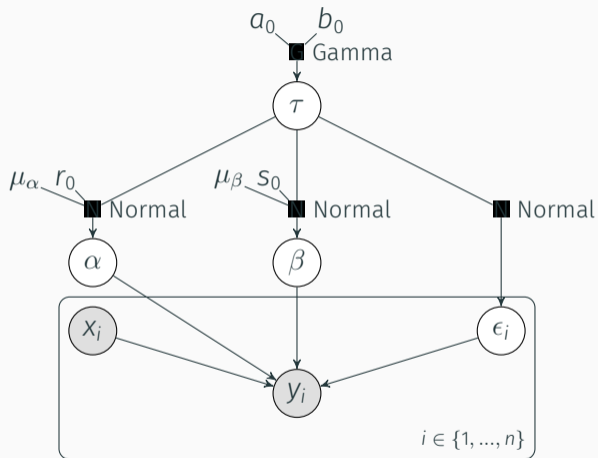
$$\alpha \sim N(\mu_\alpha, (r_0\tau)^{-1}),$$

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

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

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

# Radiata pine example



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

# Radiata pine example

Analytic evidence:

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

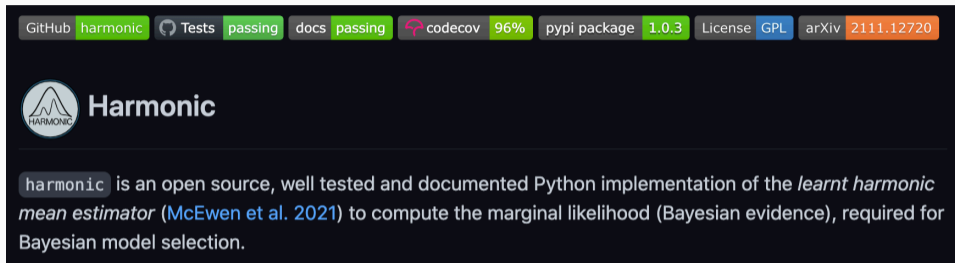
where  $\boldsymbol{\mu}_0 = (\mu_\alpha, \mu_\beta)^T$ ,  $Q_0 = \text{diag}(r_0, s_0)$ , and  $M = X^T X + Q_0$ .

# 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 $\pm 0.00072$	-301.70413 $\pm 0.00074$	8.42394 $\pm 0.00145$
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.)

# Code example

```
# Import packages
import numpy as np
import emcee
import harmonic

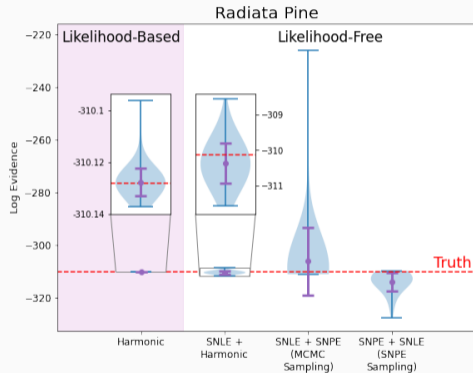
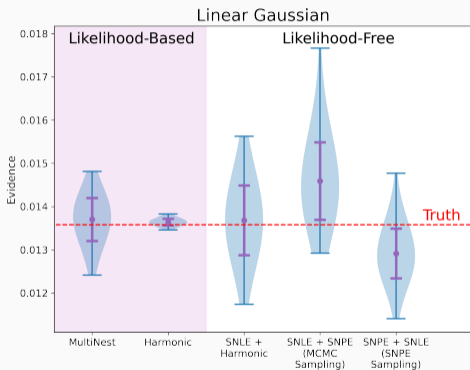
# Run MCMC sampler
sampler = emcee.EnsembleSampler(nchains, ndim, ln_posterior, args=[args])
sampler.run_mcmc(pos, samples_per_chain)
samples = np.ascontiguousarray(sampler.chain[:, nburn:, :])
lnprob = np.ascontiguousarray(sampler.lnprobability[:, nburn:])

# Set up chains
chains = harmonic.Chains(ndim)
chains.add_chains_3d(samples, lnprob)

# Fit model
chains_train, chains_test = harmonic.utils.split_data(chains, train_prop=0.05)
model = harmonic.model.KernelDensityEstimate(ndim, domain, hyper_parameters)
model.fit(chains_train.samples, chains_train.ln_posterior)

# Compute evidence
evidence = harmonic.Evidence(chains_test.nchains, model)
evidence.add_chains(chains_test)
ln_evidence, ln_evidence_std = evidence.compute_ln_evidence()
```

# Model comparison for likelihood-free inference





Proximal nested sampling for  
high-dimensional Bayesian model  
comparison

---

# Nested sampling

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  $\Omega$  into a series of **nested subspaces**.

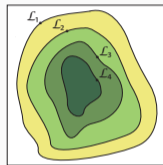
Define the prior volume  $\xi$  by  $d\xi = \pi(x)dx$ , where

$$\xi(L^*) = \int_{\Omega_{L^*}} \pi(x)dx.$$

The marginal likelihood integral can then be rewritten as

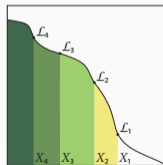
$$\mathcal{Z} = \int_0^1 \mathcal{L}(\xi)d\xi,$$

which is a **one-dimensional integral** over the prior volume  $\xi$ .



Feroz et al. (2013)

Nested subspaces



Feroz et al. (2013)

Reparameterised  
likelihood

# Constrained sampling from the prior

To compute the marginal likelihood by nested sampling, thus require strategy to generate likelihoods  $L_j$  and associated prior volumes  $\xi_j$ .

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

# Constrained sampling from the prior

To compute the marginal likelihood by nested sampling, thus require strategy to generate likelihoods  $L_j$  and associated prior volumes  $\xi_j$ .

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.

# 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  $\Omega$ .

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  $\pi$  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:

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

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

Based on **Langevin diffusion process**  $\mathcal{L}(t)$ , with  $\pi$  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**.

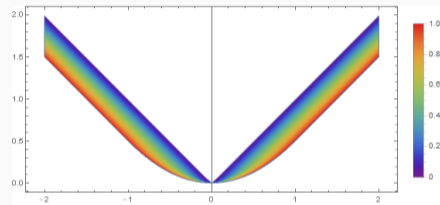
# Moreau-Yosida approximation

Moreau-Yosida approximation (envelope) of  $f$ :

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

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

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



Moreau-Yosida envelope of  $|x|$  for varying  $\lambda$  [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:

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

# 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**  $\mathcal{B}_\tau$ : term  $[x^{(k)} - \text{prox}_{\chi_{\mathcal{B}_\tau}}(x^{(k)})]$  disappears and recover usual Langevin MCMC.
2.  $x^{(k)}$  **is not in**  $\mathcal{B}_\tau$ : 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  $\mathcal{B}_\tau$ . Acts to push the Markov chain back into the constraint set  $\mathcal{B}_\tau$  if it wanders outside of it.

# Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

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

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

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

# Proximal nested sampling intuition

Recall proximal nested sampling Markov chain:

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

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

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

Many further details regarding explicit forms 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

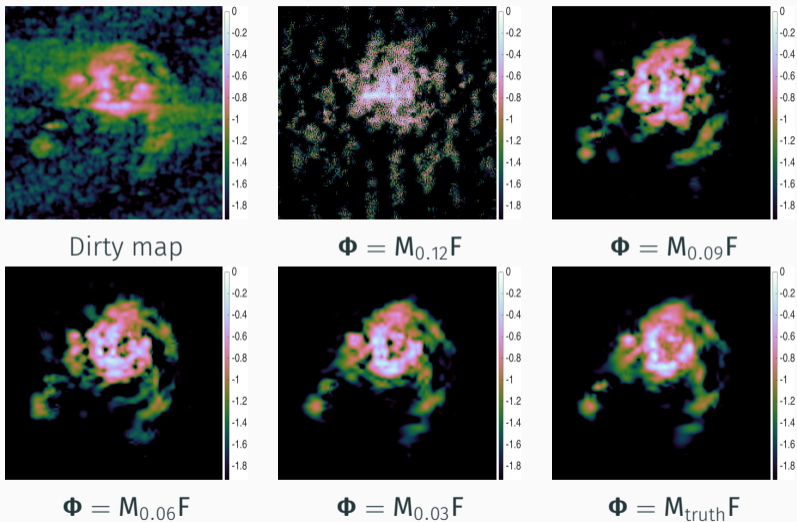
Consider ground truth model  $\Phi = \mathbf{M}_{\text{truth}}\mathbf{F}$  to simulate observational data  $\mathbf{y}$ .

However, when solving the inverse problem consider misspecified models  $\mathbf{M}_{\gamma}$ , 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

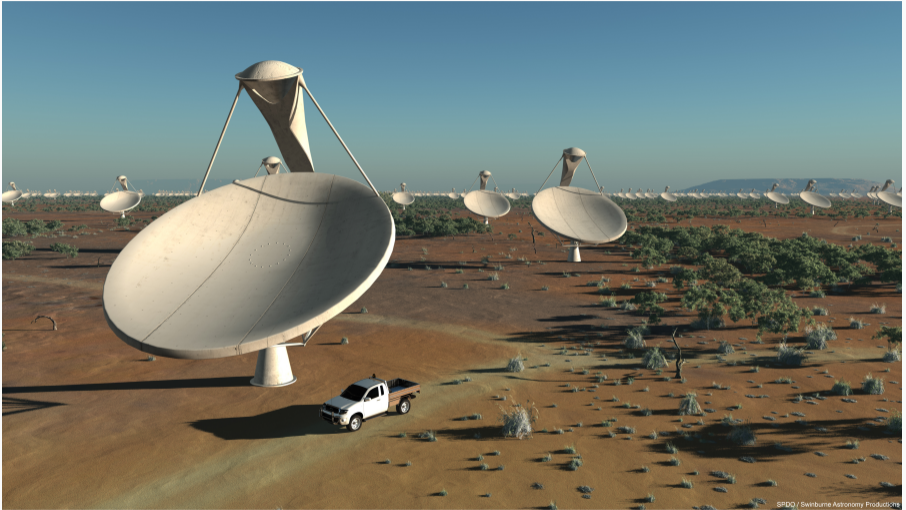
Model	$\log \mathcal{Z}$	RMSE (Requires ground truth)
$\Phi = \mathbf{M}_{\text{truth}} \mathbf{F}$	$-4.47 \times 10^3 \pm 0.08$	3.40
$\Phi = \mathbf{M}_{0.03} \mathbf{F}$	$-4.88 \times 10^3 \pm 0.08$	7.85
$\Phi = \mathbf{M}_{0.06} \mathbf{F}$	$-5.63 \times 10^3 \pm 0.08$	12.01
$\Phi = \mathbf{M}_{0.09}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

## High-dimensional Bayesian uncertainty quantification for extreme computation

---

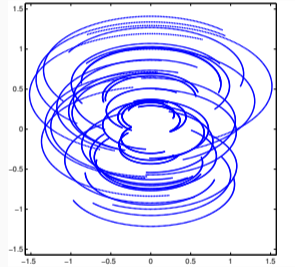
# Square Kilometre Array (SKA)



# Radio interferometric telescopes acquire “Fourier” measurements



“Fourier”  
Measurements



# Radio interferometric inverse problem

- Consider the ill-posed inverse problem of radio interferometric imaging:

$$y = \Phi x + n,$$

where  $y$  are the measured visibilities,  $\Phi$  is the linear measurement operator,  $x$  is the underlying image and  $n$  is instrumental noise.

# Radio interferometric inverse problem

- Consider the ill-posed inverse problem of radio interferometric imaging:

$$y = \Phi x + n,$$

where  $y$  are the measured visibilities,  $\Phi$  is the linear measurement operator,  $x$  is the underlying image and  $n$  is instrumental noise.

- Measurement operator, e.g.  $\Phi = GFA$ , may incorporate:
  - primary beam  $A$  of the telescope;
  - Fourier transform  $F$ ;
  - convolutional de-gridding  $G$  to interpolate to continuous  $uv$ -coordinates;
  - direction-dependent effects (DDEs)...

# Radio interferometric inverse problem

- Consider the ill-posed inverse problem of radio interferometric imaging:

$$y = \Phi x + n,$$

where  $y$  are the measured visibilities,  $\Phi$  is the linear measurement operator,  $x$  is the underlying image and  $n$  is instrumental noise.

- Measurement operator, e.g.  $\Phi = \text{GFA}$ , may incorporate:
  - primary beam  $A$  of the telescope;
  - Fourier transform  $F$ ;
  - convolutional de-gridding  $G$  to interpolate to continuous  $uv$ -coordinates;
  - direction-dependent effects (DDEs)...

**Interferometric imaging:** recover an image from noisy and incomplete Fourier measurements.



# Interferometric imaging and MAP estimation

Many interferometric imaging approaches are based on **regularisation**, *i.e.* minimising an objective function comprised of a data-fidelity penalty and a regularisation penalty.

From a Bayesian perspective this is **maximum a-posteriori (MAP) estimation**...

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

$$P(\mathbf{y}|\mathbf{x}) \propto \exp\left(-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2)\right)$$

likelihood

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

$$P(\mathbf{y}|\mathbf{x}) \propto \exp\left(-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2)\right)$$

likelihood

$$P(\mathbf{x}) \propto \exp(-R(\mathbf{x}))$$

prior

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x} | \mathbf{y}) \propto P(\mathbf{y} | \mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

$$P(\mathbf{y} | \mathbf{x}) \propto \exp\left(-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2)\right)$$

likelihood

$$P(\mathbf{x}) \propto \exp(-R(\mathbf{x}))$$

prior

Consider log-posterior:

$$\log P(\mathbf{x} | \mathbf{y}) = -\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2) - R(\mathbf{x}) + \text{const.}$$

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x} | \mathbf{y}) \propto P(\mathbf{y} | \mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

$$P(\mathbf{y} | \mathbf{x}) \propto \exp\left(-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2)\right)$$

likelihood

$$P(\mathbf{x}) \propto \exp(-R(\mathbf{x}))$$

prior

Consider log-posterior:

$$\log P(\mathbf{x} | \mathbf{y}) = -\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2) - R(\mathbf{x}) + \text{const.}$$

**MAP estimator:**

$$\mathbf{x}_{\text{map}} = \arg \max_{\mathbf{x}} \left[ \log P(\mathbf{y} | \mathbf{x}) \right]$$

# MAP estimation and regularisation

Start with Bayes Theorem (ignore normalising evidence):

$$P(\mathbf{x} | \mathbf{y}) \propto P(\mathbf{y} | \mathbf{x})P(\mathbf{x}), \quad \text{i.e. posterior} \propto \text{likelihood} \times \text{prior}$$

Define likelihood (assuming Gaussian noise) and prior:

$$P(\mathbf{y} | \mathbf{x}) \propto \exp\left(-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2)\right)$$

likelihood

$$P(\mathbf{x}) \propto \exp(-R(\mathbf{x}))$$

prior

Consider log-posterior:

$$\log P(\mathbf{x} | \mathbf{y}) = -\|\mathbf{y} - \Phi\mathbf{x}\|_2^2 / (2\sigma^2) - R(\mathbf{x}) + \text{const.}$$

**MAP estimator:**

$$\mathbf{x}_{\text{map}} = \arg \max_{\mathbf{x}} \left[ \log P(\mathbf{y} | \mathbf{x}) \right] = \arg \min_{\mathbf{x}} \left[ \underbrace{\|\mathbf{y} - \Phi\mathbf{x}\|_2^2}_{\text{data fidelity}} + \underbrace{\lambda R(\mathbf{x})}_{\text{regulariser}} \right]$$



# CLEAN and MEM as MAP estimators

- CLEAN

Consider the sparse prior:  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_0)$ .

# CLEAN and MEM as MAP estimators

- CLEAN

Consider the sparse prior:  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_0)$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{clean}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0 \right]$$

(Laplace prior  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_1)$  is good proxy)

# CLEAN and MEM as MAP estimators

- **CLEAN**

Consider the sparse prior:  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_0)$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{clean}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0 \right]$$

(Laplace prior  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_1)$  is good proxy)

- **MEM**

Consider the entropic prior:  $P(\mathbf{x}) \propto \exp(-\beta \mathbf{x}^\dagger \log \mathbf{x})$ .

# CLEAN and MEM as MAP estimators

- **CLEAN**

Consider the sparse prior:  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_0)$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{clean}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0 \right]$$

(Laplace prior  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_1)$  is good proxy)

- **MEM**

Consider the entropic prior:  $P(\mathbf{x}) \propto \exp(-\beta \mathbf{x}^\dagger \log \mathbf{x})$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{mem}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \mathbf{x}^\dagger \log \mathbf{x} \right]$$

# CLEAN and MEM as MAP estimators

- **CLEAN**

Consider the sparse prior:  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_0)$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{clean}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0 \right]$$

(Laplace prior  $P(\mathbf{x}) \propto \exp(-\beta \|\mathbf{x}\|_1)$  is good proxy)

- **MEM**

Consider the entropic prior:  $P(\mathbf{x}) \propto \exp(-\beta \mathbf{x}^\dagger \log \mathbf{x})$ .

Corresponding MAP estimator is:

$$\mathbf{x}_{\text{mem}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \mathbf{x}^\dagger \log \mathbf{x} \right]$$

(In practice some differences: CLEAN does not solve MAP problem exactly;  
MEM considered in RI imposes additional constraints.)

# Sparse regularisation (cf. compressive sensing)

Sparse **synthesis** regularisation problem:

$$x_{\text{synthesis}} = \Psi \times \arg \min_{\alpha} \left[ \|y - \Phi \Psi \alpha\|_2^2 + \lambda \|\alpha\|_1 \right]$$

synthesis framework

where consider sparsifying (e.g. wavelet) representation of image:  $x = \Psi \alpha$  .

# Sparse regularisation (cf. compressive sensing)

Sparse **synthesis** regularisation problem:

$$\mathbf{x}_{\text{synthesis}} = \Psi \times \arg \min_{\alpha} \left[ \|\mathbf{y} - \Phi \Psi \alpha\|_2^2 + \lambda \|\alpha\|_1 \right]$$

synthesis framework

where consider sparsifying (e.g. wavelet) representation of image:  $\mathbf{x} = \Psi \alpha$ .

Sparse **analysis** regularisation problem (Elad *et al.* 2007, Nam *et al.* 2012):

$$\mathbf{x}_{\text{analysis}} = \arg \min_{\mathbf{x}} \left[ \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda \|\Psi^\dagger \mathbf{x}\|_1 \right]$$

analysis framework

# Sparse regularisation (cf. compressive sensing)

Sparse **synthesis** regularisation problem:

$$x_{\text{synthesis}} = \Psi \times \arg \min_{\alpha} \left[ \|y - \Phi \Psi \alpha\|_2^2 + \lambda \|\alpha\|_1 \right]$$

synthesis framework

where consider sparsifying (e.g. wavelet) representation of image:  $x = \Psi \alpha$ .

Sparse **analysis** regularisation problem (Elad *et al.* 2007, Nam *et al.* 2012):

$$x_{\text{analysis}} = \arg \min_x \left[ \|y - \Phi x\|_2^2 + \lambda \|\Psi^\dagger x\|_1 \right]$$

analysis framework

More sophisticated extensions (e.g. overcomplete dictionaries, constrained vs unconstrained, re-weighting).



# MAP estimation vs MCMC sampling

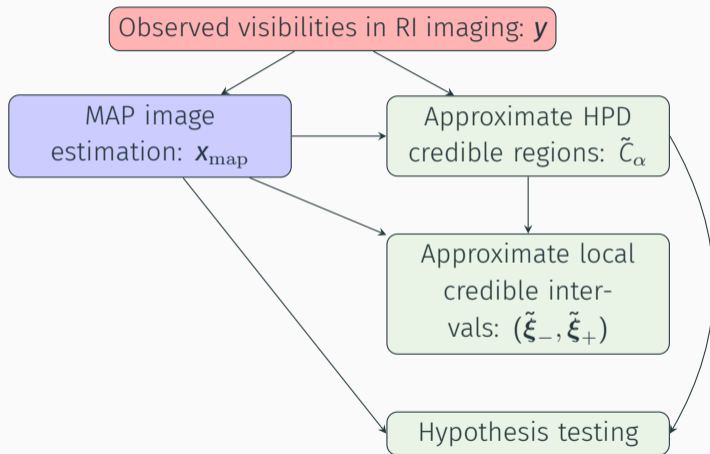
## MAP estimation

- + Based on optimization so computationally efficient.
- Does not traditionally provide uncertainties.

## MCMC sampling

- Based on sampling so computationally demanding.
- + Recover full posterior distribution.

# MAP estimation and uncertainty quantification



# Approximate Bayesian credible regions for MAP estimation

Combine **uncertainty quantification and scalable MAP estimation** (sparse regularisation) to scale to big-data (Cai, Pereyra & McEwen 2017, 2018; [arXiv:1711.04819](#); [arXiv:1811.02514](#)).

# Approximate Bayesian credible regions for MAP estimation

Combine **uncertainty quantification and scalable MAP estimation** (sparse regularisation) to scale to big-data (Cai, Pereyra & McEwen 2017, 2018; [arXiv:1711.04819](#); [arXiv:1811.02514](#)).

Let  $C_\alpha$  denote the **highest posterior density (HPD) Bayesian credible region** with confidence level  $(1 - \alpha)\%$  defined by posterior iso-contour:  $C_\alpha = \{\mathbf{x} : g(\mathbf{x}) \leq \gamma_\alpha\}$ .

# Approximate Bayesian credible regions for MAP estimation

Combine **uncertainty quantification and scalable MAP estimation** (sparse regularisation) to scale to big-data (Cai, Pereyra & McEwen 2017, 2018; [arXiv:1711.04819](#); [arXiv:1811.02514](#)).

Let  $C_\alpha$  denote the **highest posterior density (HPD) Bayesian credible region** with confidence level  $(1 - \alpha)\%$  defined by posterior iso-contour:  $C_\alpha = \{\mathbf{x} : g(\mathbf{x}) \leq \gamma_\alpha\}$ .

Analytic approximation of  $\gamma_\alpha$ :

$$\tilde{\gamma}_\alpha = g(\mathbf{x}^*) + N(\tau_\alpha + 1)$$

where  $\tau_\alpha = \sqrt{16 \log(3/\alpha)/N}$  and  $\alpha \in (4\exp(-N/3), 1)$  (Pereyra 2016b). Define **approximate HPD regions** by  $\tilde{C}_\alpha = \{\mathbf{x} : g(\mathbf{x}) \leq \tilde{\gamma}_\alpha\}$ .

**Compute  $\mathbf{x}^*$**  by MAP estimation (optimization), then **estimate local Bayesian credible intervals** and perform **uncertainty quantification** using approximate HPD regions.

# Local Bayesian credible intervals for MAP estimation

## Local Bayesian credible intervals for MAP estimation

(Cai, Pereyra & McEwen 2017, 2018; [arXiv:1711.04819](#); [arXiv:1811.02514](#))

Let  $\Omega$  define the area (or pixel) over which to compute the credible interval  $(\tilde{\xi}_-, \tilde{\xi}_+)$  and  $\zeta$  be an index vector describing  $\Omega$  (i.e.  $\zeta_i = 1$  if  $i \in \Omega$  and 0 otherwise).

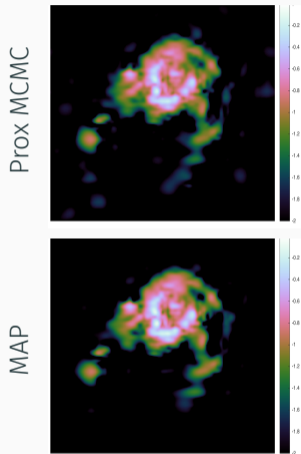
Consider the test image with the  $\Omega$  region replaced by constant value  $\xi$ :

$$\mathbf{x}' = \mathbf{x}^*(\mathcal{I} - \zeta) + \xi\zeta .$$

Given  $\tilde{\gamma}_\alpha$  and  $\mathbf{x}^*$ , compute the credible interval by

$$\begin{aligned}\tilde{\xi}_- &= \min_{\xi} \{ \xi \mid g\mathbf{y}(\mathbf{x}') \leq \tilde{\gamma}_\alpha, \forall \xi \in [-\infty, +\infty) \}, \\ \tilde{\xi}_+ &= \max_{\xi} \{ \xi \mid g\mathbf{y}(\mathbf{x}') \leq \tilde{\gamma}_\alpha, \forall \xi \in [-\infty, +\infty) \} .\end{aligned}$$

# Local credible intervals for M31 experiment



(a) point estimators

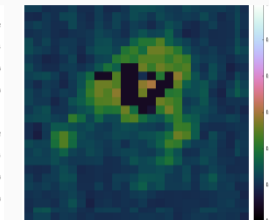
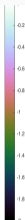
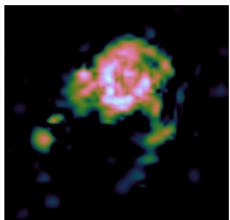
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

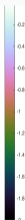
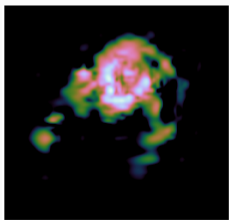
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for M31 experiment

Prox MCMC



MAP



(a) point estimators

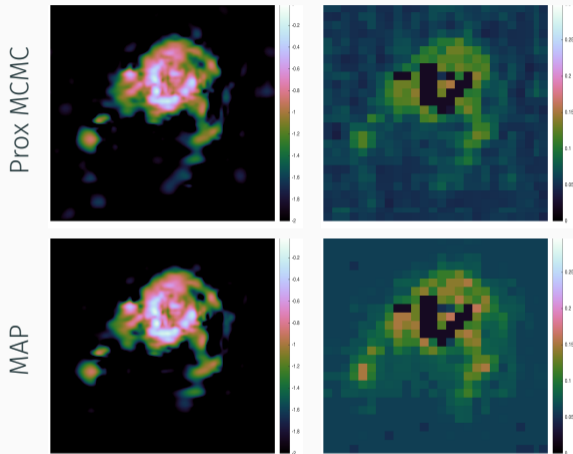
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

(d) local credible interval  
(grid size  $30 \times 30$  pixels)



# Local credible intervals for M31 experiment



(a) point estimators

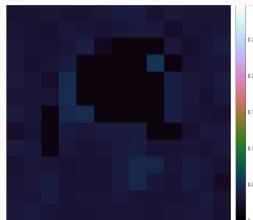
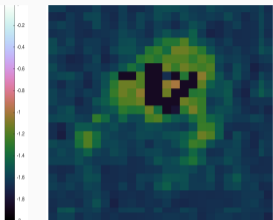
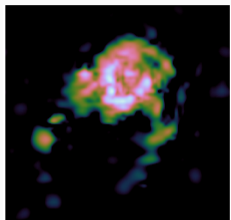
(b) local credible interval  
(grid size 10 × 10 pixels)

(c) local credible interval  
(grid size 20 × 20 pixels)

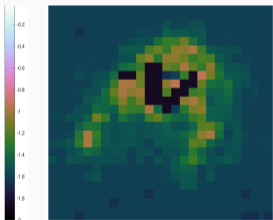
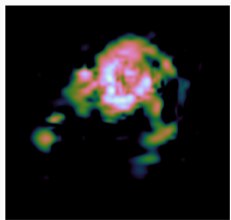
(d) local credible interval  
(grid size 30 × 30 pixels)

# Local credible intervals for M31 experiment

Prox MCMC



MAP



(a) point estimators

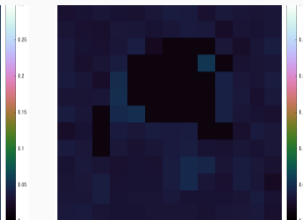
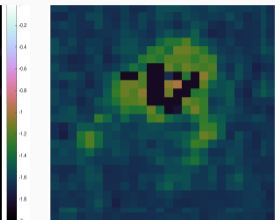
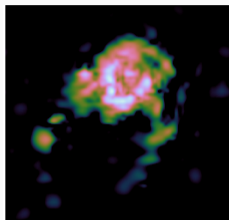
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

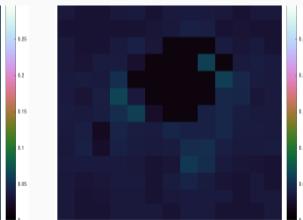
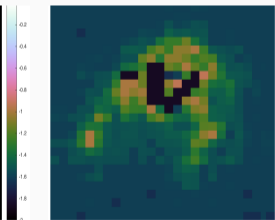
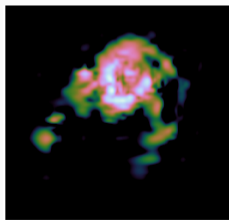
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for M31 experiment

Prox MCMC



MAP



(a) point estimators

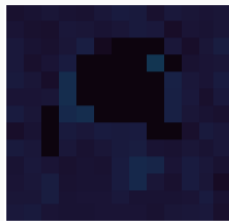
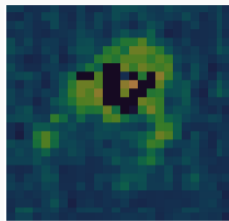
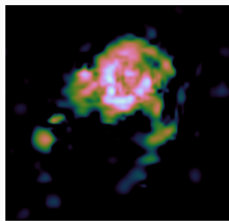
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

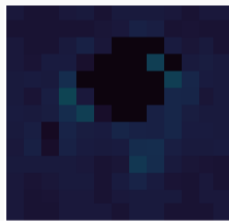
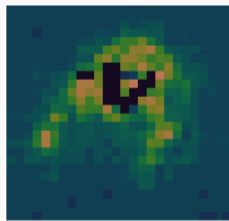
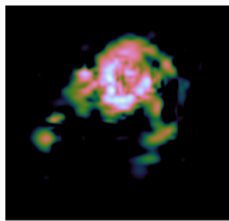
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for M31 experiment

Prox MCMC



MAP



(a) point estimators

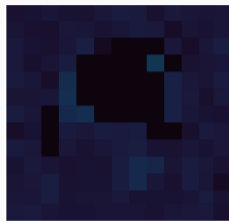
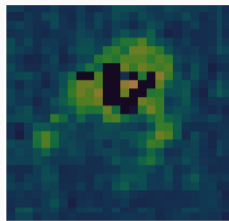
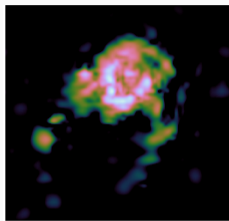
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

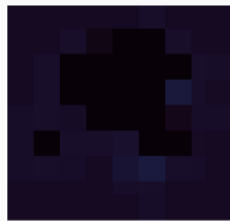
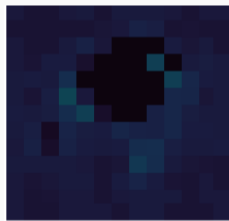
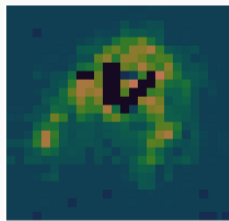
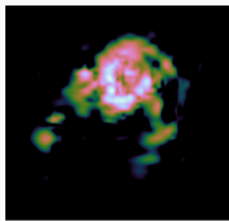
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for M31 experiment

Prox MCMC



MAP



(a) point estimators

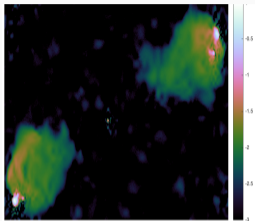
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

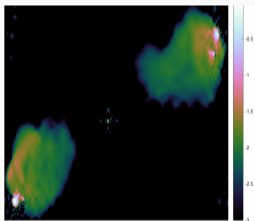
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for Cygnus A experiment

Prox MCMC



MAP



(a) point estimators

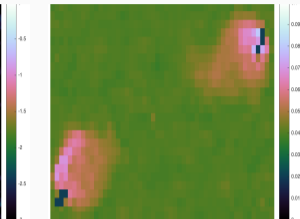
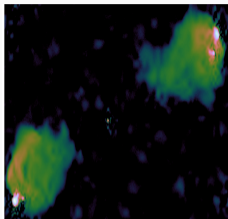
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

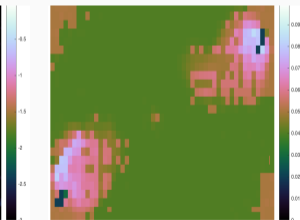
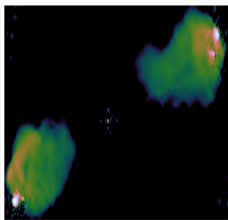
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for Cygnus A experiment

Prox MCMC



MAP



(a) point estimators

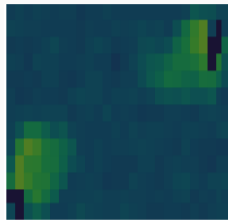
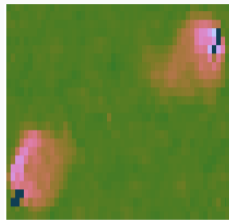
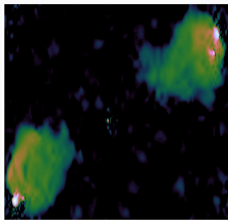
(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

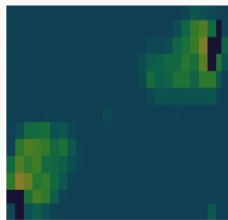
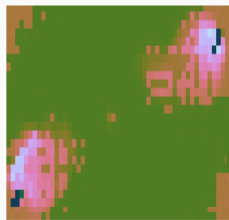
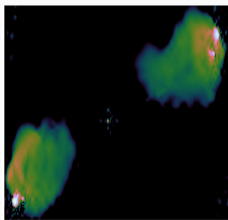
(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Local credible intervals for Cygnus A experiment

Prox MCMC



MAP



(a) point estimators

(b) local credible interval  
(grid size  $10 \times 10$  pixels)

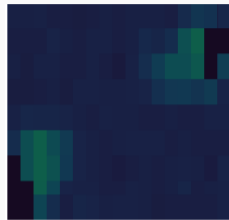
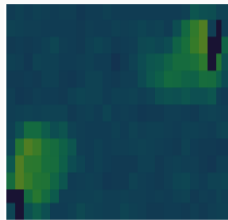
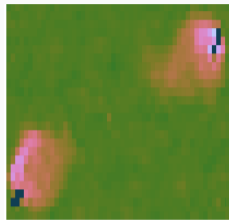
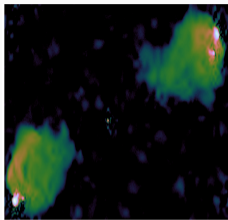
(c) local credible interval  
(grid size  $20 \times 20$  pixels)

(d) local credible interval  
(grid size  $30 \times 30$  pixels)

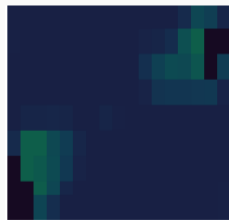
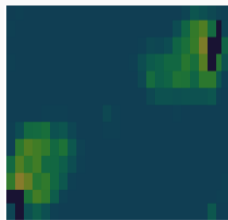
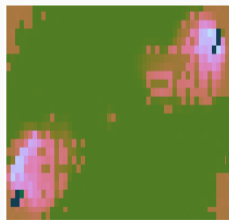
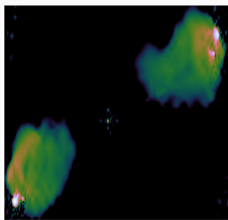


# Local credible intervals for Cygnus A experiment

Prox MCMC



MAP



(a) point estimators

(b) local credible interval  
(grid size  $10 \times 10$  pixels)

(c) local credible interval  
(grid size  $20 \times 20$  pixels)

(d) local credible interval  
(grid size  $30 \times 30$  pixels)

# Hypothesis testing

Is structure in an image physical or an artifact?

# Hypothesis testing

Is structure in an image physical or an artifact?

Perform **hypothesis tests** of image structure using Bayesian credible regions (Pereyra 2016b).

# Hypothesis testing

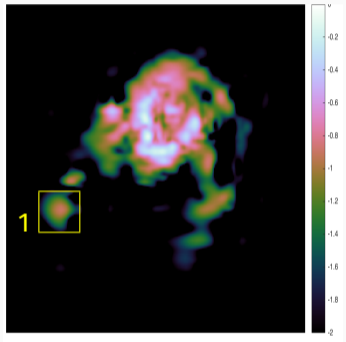
Is structure in an image physical or an artifact?

Perform **hypothesis tests** of image structure using Bayesian credible regions (Pereyra 2016b).

## Hypothesis testing of physical structure

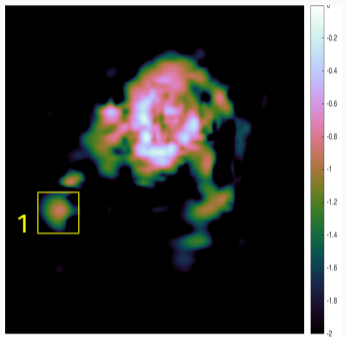
1. Remove structure of interest from recovered image  $\mathbf{x}^*$ .
2. Inpaint background (noise) into region, yielding surrogate image  $\mathbf{x}'$ .
3. Test whether  $\mathbf{x}' \in C_\alpha$ :
  - If  $\mathbf{x}' \notin C_\alpha$  then reject hypothesis that structure is an artifact with confidence  $(1 - \alpha)\%$ , *i.e.* structure most likely physical.
  - If  $\mathbf{x}' \in C_\alpha$  uncertainly too high to draw strong conclusions about the physical nature of the structure.

# Hypothesis testing for M31 experiment

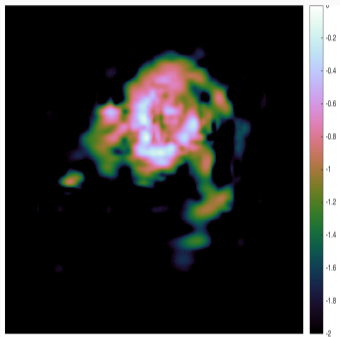


Recovered image

# Hypothesis testing for M31 experiment

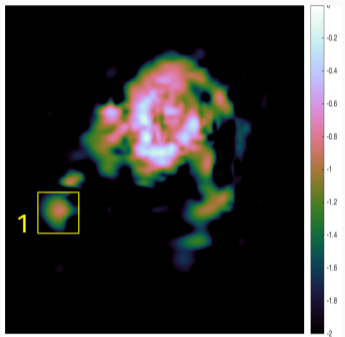


Recovered image

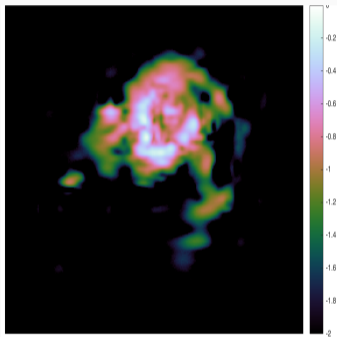


Surrogate with region removed

# Hypothesis testing for M31 experiment



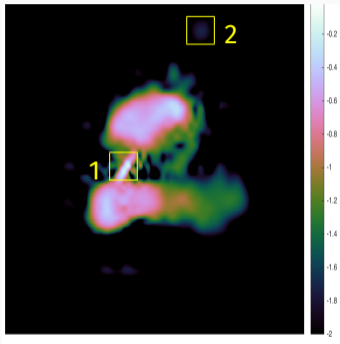
Recovered image



Surrogate with region removed

1. Reject null hypothesis  
⇒ structure physical

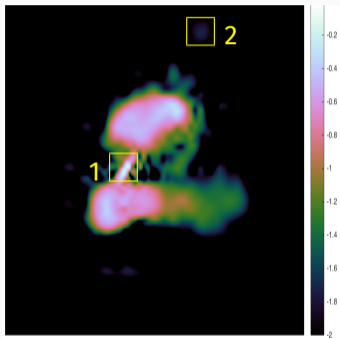
# Hypothesis testing for 3C288 experiment



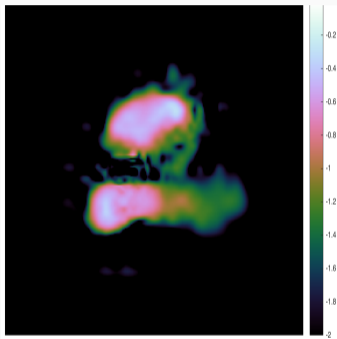
Recovered image



# Hypothesis testing for 3C288 experiment

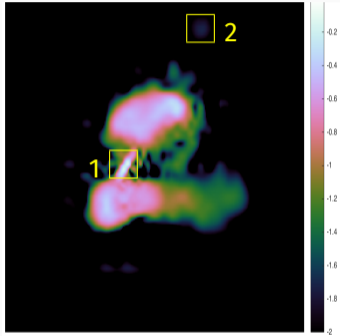


Recovered image

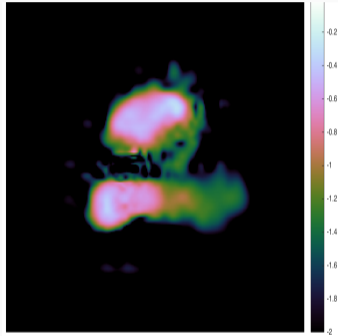


Surrogate with region removed

# Hypothesis testing for 3C288 experiment



Recovered image



Surrogate with region removed

1. Reject null hypothesis  
⇒ **structure physical**
2. Cannot reject null hypothesis  
⇒ **cannot make strong statistical statement about origin of structure**

# Computation time

CPU time in minutes for Proximal MCMC sampling and MAP estimation

Image	Method	CPU time	
		Analysis	Synthesis
Cygnus A	P-MALA	2274	1762
	MYULA	1056	942
	MAP	.07	.04
M31	P-MALA	1307	944
	MYULA	618	581
	MAP	.03	.02
3C288	P-MALA	1144	881
	MYULA	607	538
	MAP	.03	.02

# Summary

1. Learnt harmonic mean estimator for Bayesian model comparison  
(McEwen *et al.* 2021; [arXiv:2111.12720](#))
2. Proximal nested sampling for high-dimensional Bayesian model comparison  
(Cai, McEwen & Pereyra 2021; [arXiv:2106.03646](#))
3. High-dimensional Bayesian uncertainty quantification for extreme computation  
(Cai, Pereyra & McEwen 2017, 2018; [arXiv:1711.04819](#), [arXiv:1811.02514](#))