



Accelerated nested sampling with β -flows

Metha Prathaban
myp23@cam.ac.uk





About Me

- ▶ 3rd year PhD student
- ▶ Work on Bayesian numerical method development in context of GWs

Current work is in collaboration with Will Handley and Harry Bevins.

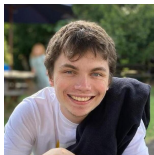




Table of Contents

Nested sampling

Accelerating NS

β -flows



Table of Contents

Nested sampling

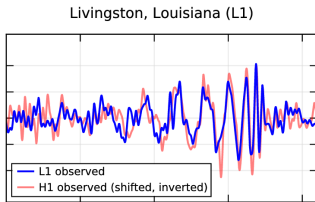
Accelerating NS

β -flows



Bayes' Theorem

Given some model \mathcal{M} and observed signal \mathcal{D} , Bayes' theorem enables us to relate the **posterior** probability of the set of parameters θ which generated the signal to the **likelihood** of the \mathcal{D} given θ and the **prior** probability of θ given \mathcal{M} :



$$\mathcal{P}(\theta|\mathcal{D}, \mathcal{M}) = \frac{P(\mathcal{D}|\theta, \mathcal{M})P(\theta|\mathcal{M})}{P(\mathcal{D}|\mathcal{M})} = \frac{\mathcal{L}(\mathcal{D}|\theta)\pi(\theta)}{\mathcal{Z}} \quad (1)$$

The **evidence**, \mathcal{Z} , plays a key role in model comparison.

```

1 {
2   "label": "lowres",
3   "outdir": "/rds/user/my23/hpc-work/acceleratedns/old_terminate/lowres_run0",
4   "sampler": "pypolychord",
5   "log_evidence": -3969.4096946457507,
6   "log_evidence_err": 0.351653537564021,
7   "log_noise_evidence": -4051.78662163637,
8   "log_bayes_factor": 82.3769269906191,
9   "priors": {
10    "mass 1": {
11      "prior": true,
12      "module": "bilby.core.prior.base",
13      "name": "Constraint",
14      "kwargs": {
15        "minimum": 5,
16        "maximum": 100,
17        "name": "mass_1",
18        "latex_label": "$m_{1s}$",
19        "unit": null
20      }
21    },
22    "mass 2": {
23      "prior": true,
24      "module": "bilby.core.prior.base",
25      "name": "Constraint",
26      "kwargs": {
27        "minimum": 5,
28        "maximum": 100,
29        "name": "mass_2",
30        "latex_label": "$m_{2s}$",

```

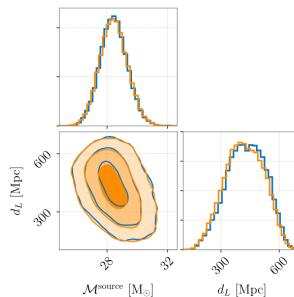


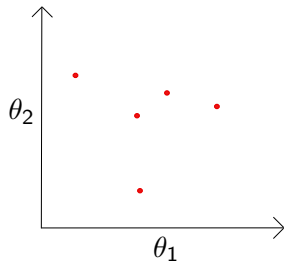
Figure 9. Posterior probability distributions for source-frame chirp mass $\mathcal{M}^{\text{source}}$ and luminosity distance d_L for GW150914. We display posteriors obtained using Bilby in orange, and LALInference posteriors in blue. We reweight the LALInference posteriors to the Bilby default priors using the procedure outlined in Appendix C. The one-dimensional JS divergence on chirp mass \mathcal{M} and luminosity distance d_L for this event are $\text{JS}_{\mathcal{M}} = 0.0017 \text{ nat}$ and $\text{JS}_{d_L} = 0.0015 \text{ nat}$.



Nested sampling (NS)

Nested sampling first and foremost calculates evidence, $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$.

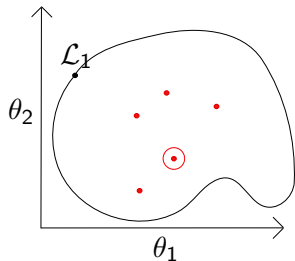
- Prior is populated with set of 'live points'.





Nested sampling (NS)

Nested sampling first and foremost calculates evidence, $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$.

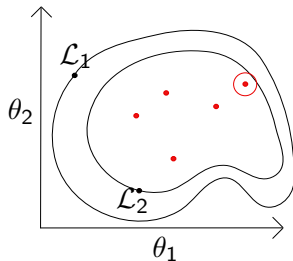


- ▶ Prior is populated with set of 'live points'.
- ▶ At each iteration i , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.



Nested sampling (NS)

Nested sampling first and foremost calculates evidence, $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$.

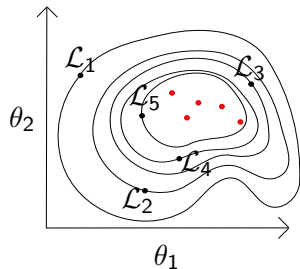


- ▶ Prior is populated with set of 'live points'.
- ▶ At each iteration i , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.



Nested sampling (NS)

Nested sampling first and foremost calculates evidence, $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$.

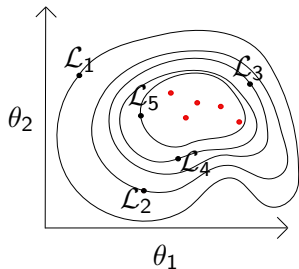


- ▶ Prior is populated with set of 'live points'.
- ▶ At each iteration i , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.
- ▶ Live points compress exponentially towards peak of likelihood.



Nested sampling (NS)

Nested sampling first and foremost calculates evidence, $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$.



- ▶ Prior is populated with set of 'live points'.
- ▶ At each iteration i , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.
- ▶ Live points compress exponentially towards peak of likelihood.
- ▶ Evidence is calculated as weighted sum over deleted ('dead') points.



Table of Contents

Nested sampling

Accelerating NS

β -flows



Time of convergence of NS:

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$



Time of convergence of NS:

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

focus of this talk



Time of convergence of NS:

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$)



Time of convergence of NS:

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

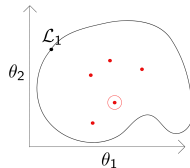
single likelihood evaluation time

compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$)



Time of convergence of NS:

drawing new live point subject to hard likelihood constraint



$$T \propto T_{\mathcal{L}} \times t_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

single likelihood evaluation time

compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$)



Time of convergence of NS:

drawing new live point subject to hard likelihood constraint

resolution

$$T \propto T_{\mathcal{L}} \times t_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

single likelihood evaluation time

compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$)



Time of convergence of NS:

drawing new live point subject to hard likelihood constraint

$$T \propto T_{\mathcal{L}} \times t_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}}$$

resolution (baked in)

single likelihood evaluation time

compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$)

(2)



Time of convergence of NS:

drawing new live point subject to hard likelihood constraint

$$T \propto T_{\mathcal{L}} \times t_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

Diagram annotations:

- Resolution (baked in) points to n_{live}
- Faster waveform models points to $T_{\mathcal{L}}$
- Compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$) points to D_{KL}



Time of convergence of NS:

$$T \propto T_{\mathcal{L}} \times t_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (2)$$

Diagram illustrating the factors affecting the time of convergence of NS (Neural Sampling):

- Resolution (baked in)**: Points to n_{live} .
- better samplers**: Points to t_{sampler} .
- faster waveform models**: Points to $T_{\mathcal{L}}$.
- compression from prior to posterior ($\approx \log \frac{V_{\pi}}{V_p}$)**: Points to D_{KL} .



Time of convergence of NS

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (3)$$

Uncertainty in $\log \mathcal{Z}$

$$\sigma \propto \sqrt{D_{\text{KL}} / n_{\text{live}}} \quad (4)$$

Precision-normalized runtime has quadratic dependence on KL divergence. 2212.01760



One way to do this (REACH):





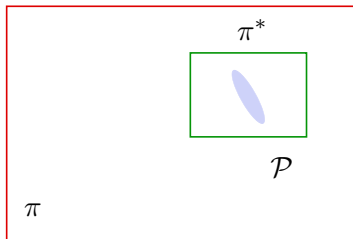
One way to do this (REACH):



- Perform low resolution (low live points) run first to roughly identify where posterior lies.



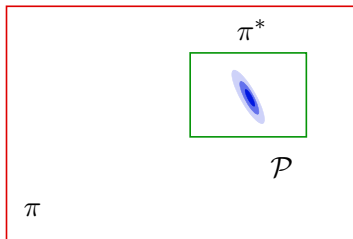
One way to do this (REACH):



- ▶ Perform low resolution (low live points) run first to roughly identify where posterior lies.
- ▶ Then set off second, high resolution, run with **narrower** box prior (much quicker).



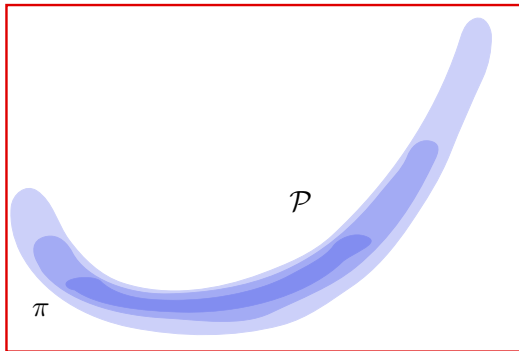
One way to do this (REACH):



- ▶ Perform low resolution (low live points) run first to roughly identify where posterior lies.
- ▶ Then set off second, high resolution, run with **narrower** box prior (much quicker).
- ▶ Evidence has **changed** (since different prior), but easy to correct (multiply new evidence by $\frac{V_{\pi^*}}{V_{\pi}}$)



When does this break down?

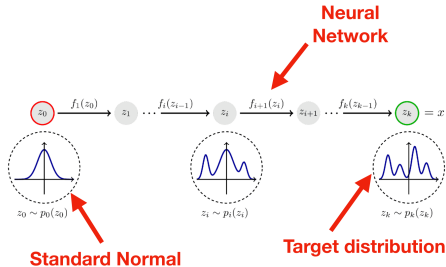


- ▶ Banana distributions, multi-modality etc.
- ▶ Precludes its use in most realistic GW cases...



- Can iterate on this by using **normalizing flows** (NF) to learn the rough posterior.

- ▶ Can iterate on this by using **normalizing flows** (NF) to learn the rough posterior.
- ▶ NFs perform density estimation, by learning a series of invertible mappings from the standard normal distribution to the target (posterior).





- ▶ Use **normalizing flows** (NF) to learn the rough posterior, and use this as our new prior, π^* .
- ▶ In this case, can't do our trick of correcting the second evidence by volume ratio, $\frac{V_{\pi^*}}{V_{\pi}}$!
- ▶ Must rely on another technique to get around this!



- ▶ Use **normalizing flows** (NF) to learn the rough posterior, and use this as our new prior, π^* .
- ▶ In this case, can't do our trick of correcting the second evidence by volume ratio, $\frac{V_{\pi^*}}{V_{\pi}}$!
- ▶ Must rely on another technique to get around this!

Posterior repartitioning (PR) can help us with this - correct likelihood accordingly to get **correct** evidence out. (see e.g. 2212.01760)

Bayesian Analysis (0000)

00, Number 0, pp. 1

Bayesian posterior repartitioning for nested sampling

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

Improving the efficiency and robustness of nested sampling using posterior repartitioning

Xi Chen · Michael Hobson · Saptarshi Das · Paul Gelderblom



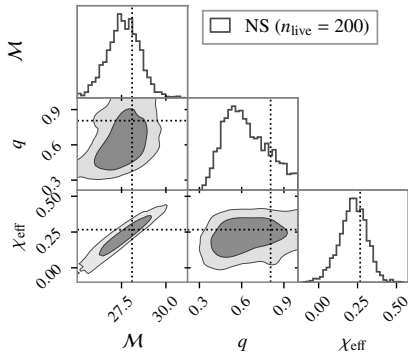
Article
SuperNest: accelerated nested sampling applied to astrophysics and cosmology[†]

Aleksandr Petrosyan^{1,2,3,*} & Will Handley^{1,2,4†}



Demo on simulated example

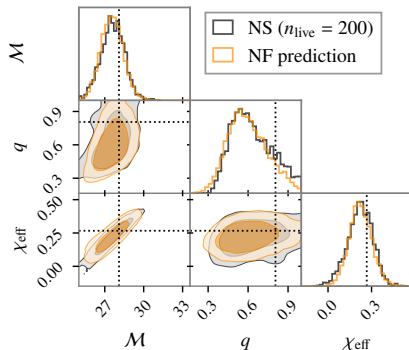
- Perform low resolution run on simulated data.





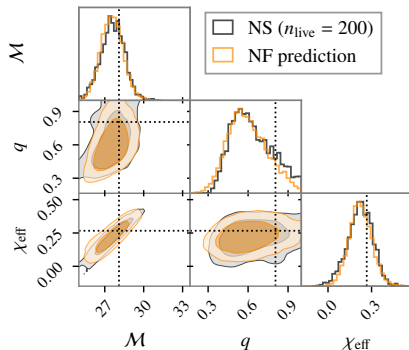
Demo on simulated example

- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.



Demo on simulated example

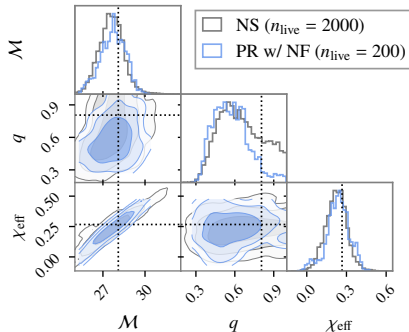
- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as 'repartitioned prior' for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).





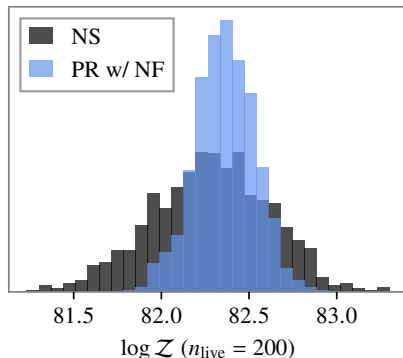
Demo on simulated example

- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as 'repartitioned prior' for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).



Demo on simulated example

- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as 'repartitioned prior' for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).

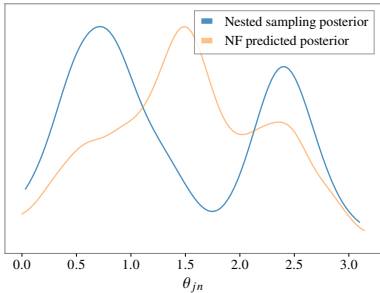


Same answer as doing a full resolution pass of NS, but **7x faster** (precision-normalized).

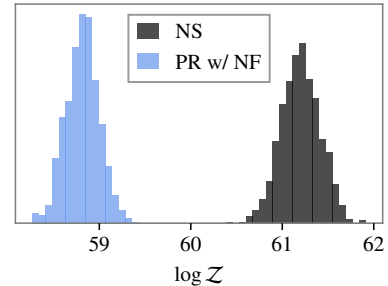
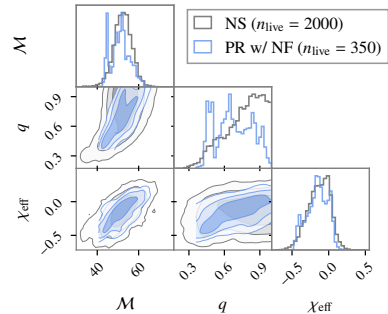


Potential pitfalls

- Still have an issue with multi-modality - if NF only learns one mode, the others can be cut off at the prior level in the high resolution run!



When the NF has been unable to properly learn the multi-modality, we can get biased posteriors and evidences (GW191222):





In order to improve the robustness of the method:

- ▶ Repartitioned prior should ideally be able to **widen itself adaptively at runtime** to mitigate missed modes and badly learned posteriors.



Table of Contents

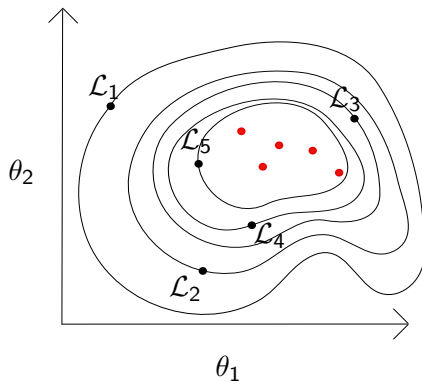
Nested sampling

Accelerating NS

β -flows



β -flows vs typical NFs

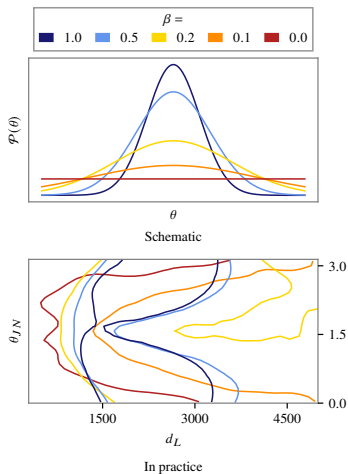


- ▶ Nested sampling sees tip to tail of the posterior in a systematic way, and NS has deep tails.
- ▶ NS can be used to train a specialized form of **conditional NFs** that can better learn these deep tail events.
- ▶ β -flows are new and only used in this work so far, though broadly applicable.



- ▶ β -flow can predict the posterior better in the tails.
- ▶ Can **widen themselves adaptively** at runtime.
- ▶ Changing a parameter β , can set repartitioned prior anywhere between posterior ($\beta = 1$) and the original prior ($\beta = 0$).

$$p(\beta) \propto \mathcal{L}^\beta \pi, \quad \beta \in [0, 1] \quad (5)$$



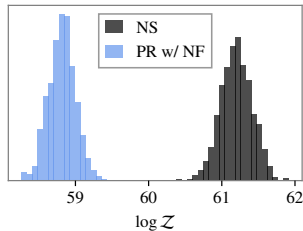
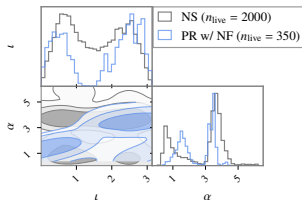
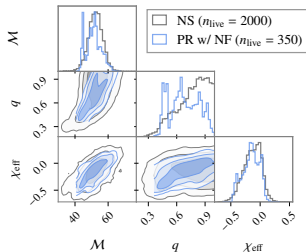


- ▶ Using β -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over β now fixes the problem.
- ▶ Although β -flow also doesn't manage to learn the full multi-modality of posterior, we can preferentially sample from lower β s (i.e. wider distributions that include missed modes).



GW191222 (again)

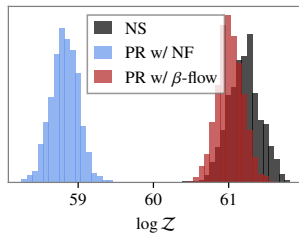
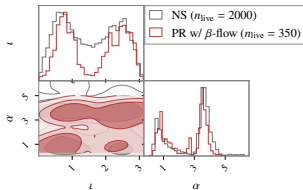
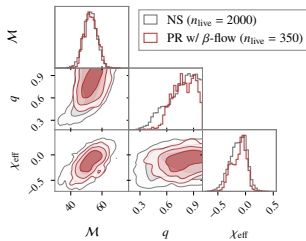
- ▶ Using β -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over β now fixes the problem.
- ▶ Although β -flow also doesn't manage to learn the full multi-modality of posterior, we can preferentially sample from lower β s (i.e. wider distributions that include missed modes).





GW191222 (again)

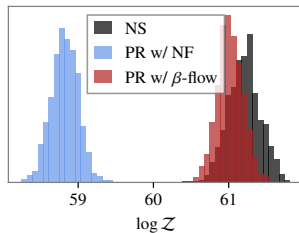
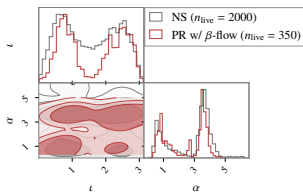
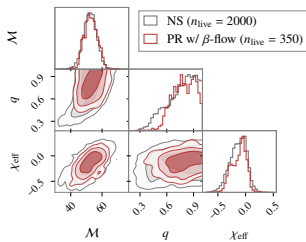
- ▶ Using β -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over β now fixes the problem.
- ▶ Although β -flow also doesn't manage to learn the full multi-modality of posterior, we can preferentially sample from lower β s (i.e. wider distributions that include missed modes).





GW191222 (again)

- ▶ Using β -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over β now fixes the problem.
- ▶ Although β -flow also doesn't manage to learn the full multi-modality of posterior, we can preferentially sample from lower β s (i.e. wider distributions that include missed modes).



Only **2x** (precision-normalized) as fast as normal NS for this real example, but robust.



- ▶ Several ways to reduce NS runtime, including reducing amount of compression from prior to posterior.
 - ▶ Can perform low resolution run to identify rough posterior, learn distribution with a flow, and perform high resolution run with this updated prior.
 - ▶ Use posterior repartitioning to get correct evidences out, despite changed prior.
 - ▶ Can achieve order of magnitude speedups on realistic GW examples.



- ▶ Several ways to reduce NS runtime, including reducing amount of compression from prior to posterior.
 - ▶ Can perform low resolution run to identify rough posterior, learn distribution with a flow, and perform high resolution run with this updated prior.
 - ▶ Use posterior repartitioning to get correct evidences out, despite changed prior.
 - ▶ Can achieve order of magnitude speedups on realistic GW examples.
- ▶ Introduced β -flows:
 - ▶ Conditional normalizing flow, trained with whole NS run
 - ▶ Better at deep tail events
 - ▶ First application is in this paper, but their use is much broader!



Thank you for listening!





Posterior repartitioning (PR)

Unlike other sampling algorithms, such as Metropolis-Hastings or Hamiltonian Monte Carlo, NS distinguishes between \mathcal{L} and π by 'sampling from the prior, subject to the hard likelihood constraint, \mathcal{L} '.

But evidence and posteriors only depend on product of \mathcal{L} and π :

$$\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta \quad (6)$$

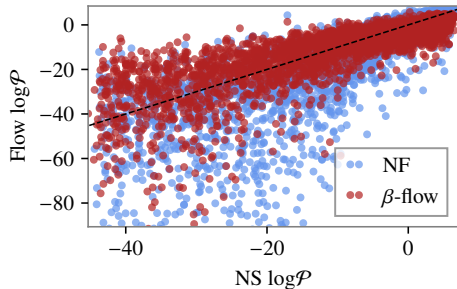
$$\mathcal{P}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{\mathcal{Z}} \quad (7)$$

Therefore, we are free to redefine the likelihood and prior however we like - as long as the product is the same! [arXiv:1908.04655](https://arxiv.org/abs/1908.04655)

$$\tilde{\mathcal{Z}} = \int \tilde{\mathcal{L}}(\theta)\tilde{\pi}(\theta)d\theta = \int \mathcal{L}(\theta)\pi(\theta)d\theta = \mathcal{Z} \quad (8)$$



Better at deep tail probabilities

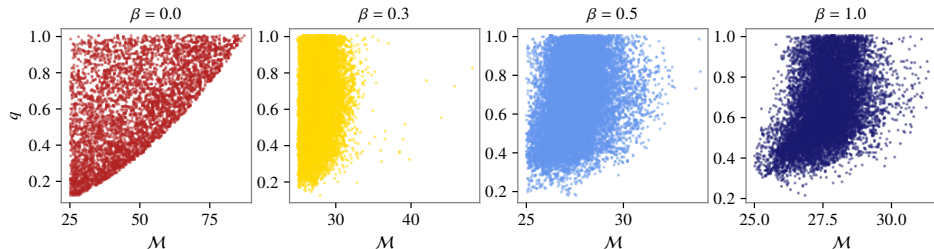


- For simulated example shown before, the β -flow is able to better predict the NS posterior probabilities.
- β -flows exhibit less scatter in the tails (low posterior probabilities) than the NFs.



What is β ?

- ▶ NS compresses step by step from prior to posterior.
- ▶ We can label these stages by a parameter β (akin to inverse temperature β in e.g. materials science).
- ▶ Sliding scale from $\beta = 0$ as the prior and $\beta = 1$ as the posterior.





Adaptive proposal

- ▶ β -flow can predict the ($\beta = 1$) posterior, similarly to NF (but better in tails).
- ▶ Can predict **any intermediate stage** of NS too.
- ▶ We **sample over** β at runtime, so proposal can now widen itself adaptively if modes have been missed!

