

---

# LINFA: A PYTHON LIBRARY FOR VARIATIONAL INFERENCE WITH NORMALIZING FLOW AND ANNEALING

---

A PREPRINT

Yu Wang, Emma R. Cobian, Jubilee Lee, Fang Liu, Jonathan D. Hauenstein, and Daniele E. Schiavazzi

Department of Applied and Computational Mathematics and Statistics  
University of Notre Dame, Notre Dame, IN, USA

July 17, 2023

## ABSTRACT

Variational inference is an increasingly popular method in statistics and machine learning for approximating probability distributions. We developed LINFA (Library for Inference with Normalizing Flow and Annealing), a Python library for variational inference to accommodate computationally expensive models and difficult-to-sample distributions with dependent parameters. We discuss the theoretical background, capabilities, and performance of LINFA in various benchmarks. LINFA is publicly available on GitHub at <https://github.com/desResLab/LINFA>.

**Keywords** Variational Inference, Normalizing Flow, Adaptive Posterior Annealing

## 1 Introduction

Generating samples from a posterior distribution is a fundamental task in Bayesian inference. The development of sampling-based algorithms from the Markov chain Monte Carlo family [1, 2, 3, 4] has made solving Bayesian inverse problems accessible to a wide audience of both researchers and practitioners. However, the number of samples required by these approaches is typically significant and the convergence of Markov chains to their stationary distribution can be slow especially in high-dimensions. Additionally, satisfactory convergence may not be always easy to quantify, even if a number of metrics have been proposed in the literature over the years. More recent paradigms have been proposed in the context of variational inference [5], where an optimization problem is formulated to determine the optimal member of a parametric family of distributions that can approximate a target posterior density. In addition, flexible approaches to parametrize variational distributions through a composition of transformations (closely related to the concept of *transport maps*, see, e.g., [6]) have reached popularity under the name of *normalizing flows* [7, 8, 9, 10, 11]. The combination of variational inference and normalizing flow has received significant recent interest in the context of general algorithms for solving inverse problems [7].

However, cases where the computational cost of evaluating the underlying probability distribution is significant occur quite often in engineering and applied sciences, for example when such evaluation requires the solution of an ordinary or partial differential equation. In such cases, inference can easily become intractable. Additionally, strong and nonlinear dependence between model parameters may result in difficult-to-sample posterior distributions characterized by features at multiple scales or by multiple modes. The LINFA library is specifically designed for cases where the model evaluation is computationally expensive. In such cases, the construction of an adaptively trained surrogate model is key to reducing the computational cost of inference [12]. In addition, LINFA provides an adaptive annealing scheduler, where temperature increments are automatically determined based on the available variational approximant of the posterior distribution. Thus, adaptive annealing makes it easier to sample from complicated densities [13].

This paper is organized as follows. The main features of the LINFA library are discussed in Section 2, followed by a brief outline of a few selected numerical tests in Section 3. Conclusions and future work are finally discussed in Section 4. The paper is completed by a brief description of the background theory and reference to the relevant

papers in Appendix A, a detailed presentation of a four benchmarks in Appendix B, and a list of all the relevant hyperparameters in Appendix C.

## 2 Capabilities

LINFA is designed as a general inference engine and allows the user to define custom input transformations, computational models, surrogates, and likelihood functions.

**1 - User-defined input parameter transformations** - Input transformations may reduce the complexity of inference and surrogate model construction in situations where the ranges of the input variables differ substantially or when the input parameters are bounded. A number of pre-defined univariate transformations are provided, i.e. *identity*, *tanh*, *linear*, and *exp*. These transformations are independently defined for each input variable, using four parameters  $(a, b, c, d)$ , providing a nonlinear transformation between the *normalized* interval  $[a, b]$  and the *physical* interval  $[c, d]$ . Additional transformations can be defined by implementing the following member functions.

- *forward* - It evaluates the transformation from the normalized to the physical space. One transformation needs to be defined for each input. For example, the list of lists

```
trsf_info = [['tanh', -7.0, 7.0, 100.0, 1500.0],
             ['tanh', -7.0, 7.0, 100.0, 1500.0],
             ['exp', -7.0, 7.0, 1.0e-5, 1.0e-2]]
```

defines a hyperbolic tangent transformation for the first two variables and an exponential transformation for the third.

- *compute\_log\_jacob\_func* - This is the log Jacobian of the transformation that needs to be included in the computation of the log posterior density to account for the additional change in volume.

**2 - User-defined computational models** - LINFA can accommodate any type of models from analytically defined posteriors with the gradient computed through automatic differentiation to legacy computational solvers for which the solution gradient is not available nor easy to compute. New models are created by implementing the methods below.

- *genDataFile* - This is a pre-processing function used to generate synthetic observations. It computes the model output corresponding to the default parameter values (usually defined as part of the model) and adds noise with a user-specified distribution. Observations will be stored in a file and are typically assigned to `model.data` so they are available for computing the log posterior.
- *solve\_t* - This function solves the model for multiple values of the *physical* input parameters specified in a matrix format (with one sample for each row and one column for each input parameter dimension).

**3 - User-defined surrogate models** - For computational models that are too expensive for online inference, LINFA provides functionalities to create, train, and fine-tune a *surrogate model*. The *Surrogate* class implements the following functionalities:

- A new surrogate model can be created using the *Surrogate* constructor.
- The *limits* (i.e. upper and lower bounds) are stored as a list of lists using the format  $[[low_0, high_0], [low_1, high_1], \dots]$ .
- A *pre-grid* is defined as an a priori selected point cloud created inside the hyper-rectangle defined by *limits*. The pre-grid can be either of type *'tensor'* (tensor product grid) where the grid order (number of points in each dimension) is defined through the argument *gridnum*, or of type *'sobol'*, in which case the variable *gridnum* defines the total number of samples.
- Surrogate model Input/Output. The two functions *surrogate\_save()* and *surrogate\_load()* are provided to save a snapshot of a given surrogate or to read it from a file.
- The *pre\_train()* function is provided to perform an initial training of the surrogate model on the pre-grid. In addition, the *update()* function is also available to re-train the model once additional training examples are available.
- The *forward()* function evaluates the surrogate model at multiple input realizations. If a transformation is defined, the surrogate should always be specified in the *normalized domain* with limits coincident with the normalized intervals.

**4 - User-defined likelihood** - A user-defined likelihood function can be defined by passing the parameters, the model, the surrogate and a coordinate transformation using *log\_density(x, model, surr, transf)* and then assigning it as a member function of the *experiment* class using:

```
exp.model_logdensity = lambda x: log_density(x, model, surr, transf).
```

- 5 - **Linear and adaptive annealing schedulers** - LINFA provides two annealing schedulers by default. The first is the '*Linear*' scheduler with constant increments. The second is the '*AdaAnn*' adaptive scheduler [13] with hyperparameters reported in Table 7. For the AdaAnn scheduler, the user can also specify a different number of parameter updates to be performed at the initial temperature  $t_0$ , final temperature  $t_1$ , and for any temperature  $t_0 < t < 1$ . Finally, the batch size (number of samples used to evaluate the expectations in the loss function) can also be differentiated for  $t = 1$  and  $t < 1$ .
- 6 - **User-defined hyperparameters** - A complete list of hyperparameters with a description of their functionality can be found in Appendix C.

### 3 Numerical benchmarks

We tested LINFA on multiple problems. These include inference on unimodal and multi-modal posterior distributions specified in closed form, ordinary differential models and dynamical systems with gradients directly computed through automatic differentiation in PyTorch, identifiable and non-identifiable physics-based models with fixed and adaptive surrogates, and high-dimensional statistical models. Some of the above tests are included with the library and systematically tested when pushing the master branch on GitHub. A detailed discussion of these test cases is provided in Appendix B. LINFA can be installed through the Python Package Index (PyPI) typing

```
pip install linfa-vi
```

To run the tests type

```
python -m unittest linfa.linfa_test_suite.NAME_example
```

where *NAME* is the name of the test case, either *trivial*, *highdim*, *rc*, *rcr*, or *adaann*.

### 4 Conclusion and Future Work

In this paper, we have introduced the LINFA library for variational inference, briefly discussed the relevant background, its capabilities, and report its performance on a number of test cases. Some interesting directions for future work are mentioned below.

Future versions will support user-defined privacy-preserving synthetic data generation and variational inference through differentially private gradient descent algorithms. This will allow the user to perform inference tasks while preserving a pre-defined privacy budget, as discussed in [14]. LINFA will also be extended to handle multiple models. This will open new possibilities to solve inverse problems combining variational inference and multi-fidelity surrogates [see, e.g., 15]. In addition, for inverse problems with significant dependence among the parameters, it is often possible to simplify the inference task by operating on manifolds of reduced dimensionality [16]. New modules for dimensionality reduction will be developed and integrated with the LINFA library. Finally, the ELBO loss typically used in variational inference has known limitations, some of which are related to its close connection with the KL divergence. Future versions of LINFA will provide the option to use alternative losses.

### Acknowledgements

The authors gratefully acknowledge the support from the NSF Big Data Science & Engineering grant #1918692 and the computational resources provided through the Center for Research Computing at the University of Notre Dame. DES also acknowledges support from NSF CAREER grant #1942662.

## A Background theory

### A.1 Variational inference with normalizing flow

Consider the problem of estimating (in a Bayesian sense) the parameters  $\mathbf{z} \in \mathcal{Z}$  of a physics-based or statistical model

$$\mathbf{x} = \mathbf{f}(\mathbf{z}) + \varepsilon,$$

from the observations  $\mathbf{x} \in \mathcal{X}$  and a known statistical characterization of the error  $\varepsilon$ . We tackle this problem with variational inference and normalizing flow. A normalizing flow (NF) is a nonlinear transformation  $F : \mathbb{R}^d \times \Lambda \rightarrow \mathbb{R}^d$

designed to map an easy-to-sample *base* distribution  $q_0(\mathbf{z}_0)$  into a close approximation  $q_K(\mathbf{z}_K)$  of a desired target posterior density  $p(\mathbf{z}|\mathbf{x})$ . This transformation can be determined by composing  $K$  bijections

$$\mathbf{z}_K = F(\mathbf{z}_0) = F_K \circ F_{K-1} \circ \dots \circ F_k \circ \dots \circ F_1(\mathbf{z}_0), \quad (1)$$

and evaluating the transformed density through the change of variable formula [see 6].

In the context of variational inference, we seek to determine an *optimal* set of parameters  $\boldsymbol{\lambda} \in \Lambda$  so that  $q_K(\mathbf{z}_K) \approx p(\mathbf{z}|\mathbf{x})$ . Given observations  $\mathbf{x} \in \mathcal{X}$ , a likelihood function  $l_{\mathbf{z}}(\mathbf{x})$  (informed by the distribution of the error  $\boldsymbol{\varepsilon}$ ) and prior  $p(\mathbf{z})$ , a NF-based approximation  $q_K(\mathbf{z})$  of the posterior distribution  $p(\mathbf{z}|\mathbf{x})$  can be computed by maximizing the lower bound to the log marginal likelihood  $\log p(\mathbf{x})$  (the so-called *evidence lower bound* or ELBO), or, equivalently, by minimizing a *free energy bound* [see, e.g., 7].

$$\begin{aligned} \mathcal{F}(\mathbf{x}) &= \mathbb{E}_{q_K(\mathbf{z}_K)} [\log q_K(\mathbf{z}_K) - \log p(\mathbf{x}, \mathbf{z}_K)] \\ &= \mathbb{E}_{q_0(\mathbf{z}_0)} [\log q_0(\mathbf{z}_0)] - \mathbb{E}_{q_0(\mathbf{z}_0)} [\log p(\mathbf{x}, \mathbf{z}_K)] - \mathbb{E}_{q_0(\mathbf{z}_0)} \left[ \sum_{k=1}^K \log \left| \det \frac{\partial \mathbf{z}_k}{\partial \mathbf{z}_{k-1}} \right| \right]. \end{aligned} \quad (2)$$

For computational convenience, normalizing flow transformations are selected to be easily invertible and their Jacobian determinant can be computed with a cost that grows linearly with the problem dimensionality. Approaches in the literature include RealNVP [8], GLOW [17], and autoregressive transformations such as MAF [18] and IAF [9].

## A.2 MAF and RealNVP

LINFA implements two widely used normalizing flow formulations, MAF [18] and RealNVP [8]. MAF belongs to the class of *autoregressive* normalizing flows. Given the latent variable  $\mathbf{z} = (z_1, z_2, \dots, z_d)$ , it assumes  $p(z_i|z_1, \dots, z_{i-1}) = \phi[(z_i - \mu_i)/e^{\alpha_i}]$ , where  $\phi$  is the standard normal distribution,  $\mu_i = f_{\mu_i}(z_1, \dots, z_{i-1})$ ,  $\alpha_i = f_{\alpha_i}(z_1, \dots, z_{i-1})$ ,  $i = 1, 2, \dots, d$ , and  $f_{\mu_i}$  and  $f_{\alpha_i}$  are masked autoencoder neural networks [MADE, 19]. In a MADE autoencoder the network connectivities are multiplied by Boolean masks so the input-output relation maintains a lower triangular structure, making the computation of the Jacobian determinant particularly simple. MAF transformations are then composed of multiple MADE layers, possibly interleaved by batch normalization layers [20], typically used to add stability during training and increase network accuracy [18].

RealNVP is another widely used flow where, at each layer the first  $d'$  variables are left unaltered while the remaining  $d - d'$  are subject to an affine transformation of the form  $\widehat{\mathbf{z}}_{d'+1:d} = \mathbf{z}_{d'+1:d} \odot e^{\boldsymbol{\alpha}} + \boldsymbol{\mu}$ , where  $\boldsymbol{\mu} = f_{\boldsymbol{\mu}}(\mathbf{z}_{1:d'})$  and  $\boldsymbol{\alpha} = f_{\boldsymbol{\alpha}}(\mathbf{z}_{d'+1:d})$  are MADE autoencoders. In this context, MAF could be seen as a generalization of RealNVP by setting  $\mu_i = \alpha_i = 0$  for  $i \leq d'$  [18].

## A.3 Normalizing flow with adaptive surrogate (NoFAS)

LINFA is designed to accommodate black-box models  $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$  between the random inputs  $\mathbf{z} = (z_1, z_2, \dots, z_d)^T \in \mathcal{Z}$  and the outputs  $(x_1, x_2, \dots, x_m)^T \in \mathcal{X}$ , and assumes  $n$  observations  $\mathbf{x} = \{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{X}$  to be available. Our goal is to infer  $\mathbf{z}$  and to quantify its uncertainty given  $\mathbf{x}$ . We employ a variational Bayesian paradigm and sample from the posterior distribution  $p(\mathbf{z}|\mathbf{x}) \propto l_{\mathbf{z}}(\mathbf{x}, \mathbf{f}) p(\mathbf{z})$ , with prior  $p(\mathbf{z})$  via normalizing flows.

This requires the evaluation of the gradient of the ELBO (2) with respect to the NF parameters  $\boldsymbol{\lambda}$ , replacing  $p(\mathbf{x}, \mathbf{z}_K)$  with  $p(\mathbf{x}|\mathbf{z}_K) p(\mathbf{z}) = l_{\mathbf{z}_K}(\mathbf{x}, \mathbf{f}) p(\mathbf{z})$ , and approximating the expectations with their MC estimates. However, the likelihood function needs to be evaluated at every MC realization, which can be costly if the model  $\mathbf{f}(\mathbf{z})$  is computationally expensive. In addition, automatic differentiation through a legacy (e.g. physics-based) solver may be an impractical, time-consuming, or require the development of an adjoint solver.

Our solution is to replace the model  $\mathbf{f}$  with a computationally inexpensive surrogate  $\widehat{\mathbf{f}} : \mathcal{Z} \times \mathcal{W} \rightarrow \mathcal{X}$  parameterized by the weights  $\mathbf{w} \in \mathcal{W}$ , whose derivatives can be obtained at a relatively low computational cost, but intrinsic bias in the selected surrogate formulation, a limited number of training examples, and locally optimal  $\mathbf{w}$  can compromise the accuracy of  $\widehat{\mathbf{f}}$ .

To resolve these issues, LINFA implements NoFAS, which updates the surrogate model adaptively by smartly weighting the samples of  $\mathbf{z}$  from NF thanks to a *memory-aware* loss function. Once a newly updated surrogate is obtained, the likelihood function is updated, leading to a new posterior distribution that will be approximated by VI-NF, producing, in turn, new samples for the next surrogate model update, and so on. Additional details can be found in [12].

#### A.4 Adaptive Annealing

Annealing is a technique to parametrically smooth a target density to improve sampling efficiency and accuracy during inference. In the discrete case, this is achieved by incrementing an *inverse temperature*  $t_k$  and setting  $p_k(\mathbf{z}, \mathbf{x}) = p^{t_k}(\mathbf{z}, \mathbf{x})$ , for  $k = 0, \dots, K$ , where  $0 < t_0 < \dots < t_K \leq 1$ . The result of exponentiation produces a smooth unimodal distribution for a sufficiently small  $t_0$ , recovering the target density as  $t_k$  approaches 1. In other words, annealing provides a continuous deformation from an easier to approximate unimodal distribution to a desired target density.

A linear annealing scheduler [see, e.g., 7] with fixed temperature increments is often used in practice, where  $t_j = t_0 + j(1 - t_0)/K$  for  $j = 0, \dots, K$  with constant increments  $\epsilon = (1 - t_0)/K$ . Intuitively, small temperature changes are desirable to carefully explore the parameter spaces at the beginning of the annealing process, whereas larger changes can be taken as  $t_k$  increases, after annealing has helped to capture important features of the target distribution (e.g., locating all the relevant modes).

The AdaAnn scheduler determines the increment  $\epsilon_k$  that approximately produces a pre-defined change in the KL divergence between two distributions annealed at  $t_k$  and  $t_{k+1} = t_k + \epsilon_k$ , respectively. Letting the KL divergence equal a constant  $\tau^2/2$ , where  $\tau$  is referred to as the *KL tolerance*, the step size  $\epsilon_k$  becomes

$$\epsilon_k = \tau / \sqrt{\mathbb{V}_{p^{t_k}}[\log p(\mathbf{z}, \mathbf{x})]}. \quad (3)$$

The denominator is large when the support of the annealed distribution  $p^{t_k}(\mathbf{z}, \mathbf{x})$  is wider than the support of the target  $p(\mathbf{z}, \mathbf{x})$ , and progressively reduces with increasing  $t_k$ . Further detail on the derivation of the expression for  $\epsilon_k$  can be found in [13].

## B Detailed numerical benchmarks

### B.1 Simple two-dimensional map with Gaussian likelihood

A model  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is chosen in this experiment having the closed-form expression

$$f(\mathbf{z}) = f(z_1, z_2) = (z_1^3/10 + \exp(z_2/3), z_1^3/10 - \exp(z_2/3))^T. \quad (4)$$

Observations  $\mathbf{x}$  are generated as

$$\mathbf{x} = \mathbf{x}^* + 0.05 |\mathbf{x}^*| \odot \mathbf{x}_0, \quad (5)$$

where  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_2)$  and  $\odot$  is the Hadamard product. We set the *true* model parameters at  $\mathbf{z}^* = (3, 5)^T$ , with output  $\mathbf{x}^* = f(\mathbf{z}^*) = (7.99, -2.59)^T$ , and simulate 50 sets of observations from (5). The likelihood of  $\mathbf{z}$  given  $\mathbf{x}$  is assumed Gaussian and we adopt a noninformative uniform prior  $p(\mathbf{z})$ . We allocate a budget of  $4 \times 4 = 16$  model solutions to the pre-grid and use the rest to adaptively calibrate  $\hat{f}$  using 2 samples every 1000 normalizing flow iterations.

Results in terms of loss profile, variational approximation, and posterior predictive distribution are shown in Figure 1.

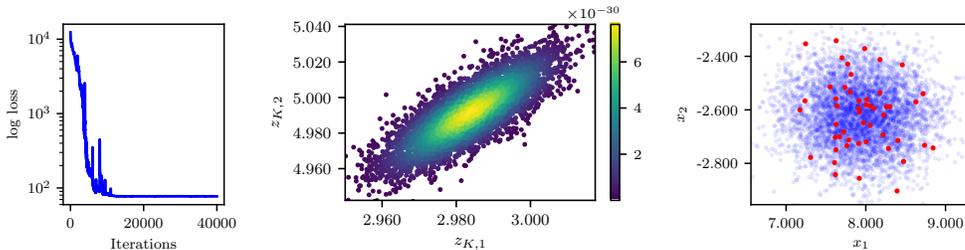


Figure 1: Results from the trivial model. Loss profile (left), posterior samples (center), and posterior predictive distribution (right).

### B.2 High-dimensional example

We consider a map  $f : \mathbb{R}^5 \rightarrow \mathbb{R}^4$  expressed as

$$f(\mathbf{z}) = \mathbf{A} \mathbf{g}(e^{\mathbf{z}}), \quad (6)$$

where  $g_i(\mathbf{r}) = (2 \cdot |2a_i - 1| + r_i)/(1 + r_i)$  with  $r_i > 0$  for  $i = 1, \dots, 5$  is the *Sobol* function [21] and  $\mathbf{A}$  is a  $4 \times 5$  matrix. We also set

$$\mathbf{a} = (0.084, 0.229, 0.913, 0.152, 0.826)^T \text{ and } \mathbf{A} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{pmatrix}.$$

The true parameter vector is set at  $\mathbf{z}^* = (2.75, -1.5, 0.25, -2.5, 1.75)^T$ . While the Sobol function is bijective and analytic,  $f$  is over-parameterized and non identifiable. This is also confirmed by the fact that the curve segment  $\gamma(t) = g^{-1}(g(\mathbf{z}^*) + \mathbf{v}t) \in Z$  gives the same model solution as  $\mathbf{x}^* = f(\mathbf{z}^*) = f(\gamma(t)) \approx (1.4910, 1.6650, 1.8715, 1.7011)^T$  for  $t \in (-0.0153, 0.0686]$ , where  $\mathbf{v} = (1, -1, 1, -1, 1)^T$ . This is consistent with the one-dimensional null-space of the matrix  $\mathbf{A}$ . We also generate synthetic observations from the Gaussian distribution

$$\mathbf{x} = \mathbf{x}^* + 0.01 \cdot |\mathbf{x}^*| \odot \mathbf{x}_0, \text{ and } \mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_5). \quad (7)$$

Results are shown in Figure 2.

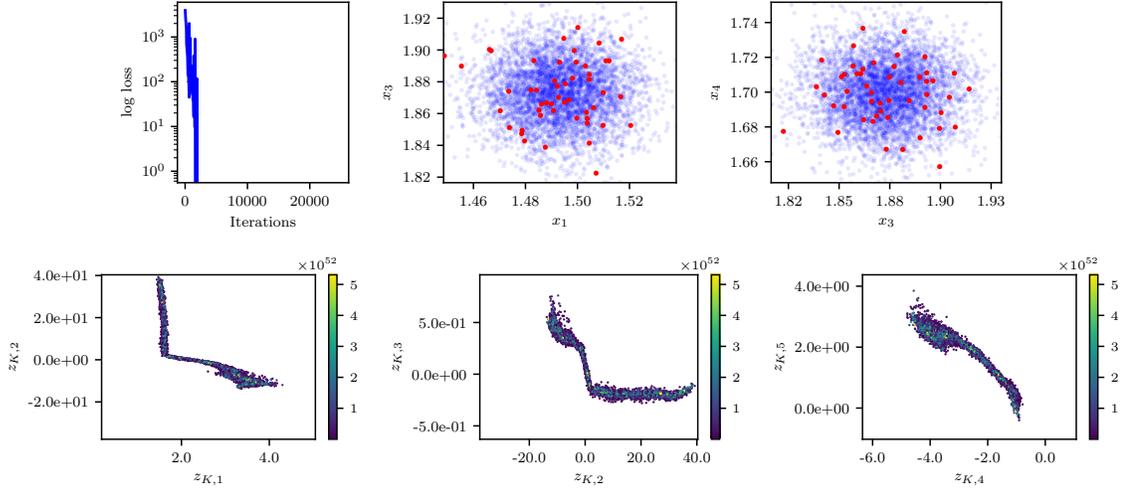


Figure 2: Results from the high-dimensional model. Loss profile, posterior samples, and posterior predictive distribution.

### B.3 Two-element Windkessel Model

The two-element Windkessel model (often referred to as the *RC* model) is the simplest representation of the human systemic circulation and requires two parameters, i.e., a resistance  $R \in [100, 1500]$  Barye·s/ml and a capacitance  $C \in [1 \times 10^{-5}, 1 \times 10^{-2}]$  ml/Barye. We provide a periodic time history of the aortic flow (see [12] for additional details) and use the *RC* model to predict the time history of the proximal pressure  $P_p(t)$ , specifically its maximum, minimum, and average values over a typical heart cycle, while assuming the distal resistance  $P_d(t)$  as a constant in time, equal to 55 mmHg. In our experiment, we set the true resistance and capacitance as  $z_{K,1}^* = R^* = 1000$  Barye·s/ml and  $z_{K,2}^* = C^* = 5 \times 10^{-5}$  ml/Barye, and determine  $P_p(t)$  from a RK4 numerical solution of the following algebraic-differential system

$$Q_d = \frac{P_p - P_d}{R}, \quad \frac{dP_p}{dt} = \frac{Q_p - Q_d}{C}, \quad (8)$$

where  $Q_p$  is the flow entering the *RC* system and  $Q_d$  is the distal flow. Synthetic observations are generated by adding Gaussian noise to the true model solution  $\mathbf{x}^* = (x_1^*, x_2^*, x_3^*) = (P_{p,\min}, P_{p,\max}, P_{p,\text{avg}}) = (78.28, 101.12, 85.75)$ , i.e.,  $\mathbf{x}$  follows a multivariate Gaussian distribution with mean  $\mathbf{x}^*$  and a diagonal covariance matrix with entries  $0.05 x_i^*$ , where  $i = 1, 2, 3$  corresponds to the maximum, minimum, and average pressures, respectively. The aim is to quantify the uncertainty in the *RC* model parameters given 50 repeated pressure measurements. We imposed a non-informative prior on  $R$  and  $C$ . Results are shown in Figure 3.

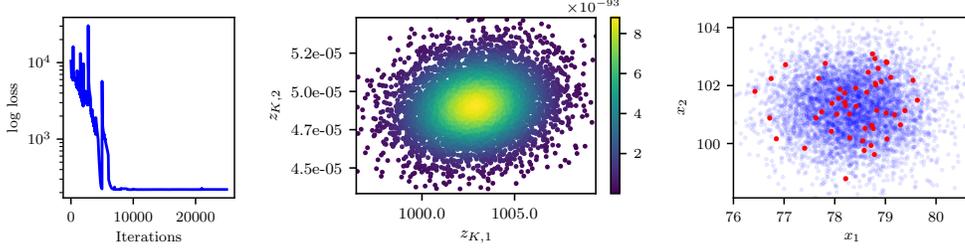


Figure 3: Results from the RC model. Loss profile (left), posterior samples (center) for R and C, and the posterior predictive distribution for  $P_{p,\min}$  and  $P_{p,\max}$  (right,  $P_{p,\text{avg}}$  not shown).

#### B.4 Three-element Windkessel Circulatory Model (NoFAS)

The three-parameter Windkessel or *RCR* model is characterized by proximal and distal resistance parameters  $R_p, R_d \in [100, 1500]$  Barye-s/ml, and one capacitance parameter  $C \in [1 \times 10^{-5}, 1 \times 10^{-2}]$  ml/Barye. This model is not identifiable. The average distal pressure is only affected by the total system resistance, i.e. the sum  $R_p + R_d$ , leading to a negative correlation between these two parameters. Thus, an increment in the proximal resistance is compensated by a reduction in the distal resistance (so the average distal pressure remains the same) which, in turn, reduces the friction encountered by the flow exiting the capacitor. An increase in the value of  $C$  is finally needed to restore the average, minimum and maximum pressure. This leads to a positive correlation between  $C$  and  $R_d$ .

The output consists of the maximum, minimum, and average values of the proximal pressure  $P_p(t)$ , i.e.,  $(P_{p,\min}, P_{p,\max}, P_{p,\text{avg}})$  over one heart cycle. The true parameters are  $z_{K,1}^* = R_p^* = 1000$  Barye-s/ml,  $z_{K,2}^* = R_d^* = 1000$  Barye-s/ml, and  $C^* = 5 \times 10^{-5}$  ml/Barye. The proximal pressure is computed from the solution of the algebraic-differential system

$$Q_p = \frac{P_p - P_c}{R_p}, \quad Q_d = \frac{P_c - P_d}{R_d}, \quad \frac{dP_c}{dt} = \frac{Q_p - Q_d}{C}, \quad (9)$$

where the distal pressure is set to  $P_d = 55$  mmHg. Synthetic observations are generated from  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = (f_1(\mathbf{z}^*), f_2(\mathbf{z}^*), f_3(\mathbf{z}^*))^T = (P_{p,\min}, P_{p,\max}, P_{p,\text{ave}})^T = (100.96, 148.02, 116.50)^T$  and  $\boldsymbol{\Sigma}$  is a diagonal matrix with entries  $(5.05, 7.40, 5.83)^T$ . The budgeted number of true model solutions is 216; the fixed surrogate model is evaluated on a  $6 \times 6 \times 6 = 216$  pre-grid while the adaptive surrogate is evaluated with a pre-grid of size  $4 \times 4 \times 4 = 64$  and the other 152 evaluations are adaptively selected.

This example also demonstrates how NoFAS can be combined with annealing for improved convergence. The results in Figure 4 are generated using the AdaAnn adaptive annealing scheduler with initial inverse temperature  $t_0 = 0.05$ , KL tolerance  $\tau = 0.01$  and a batch size of 100 samples. The number of parameter updates is set to 500, 5000 and 5 for  $t_0, t_1$  and  $t_0 < t < t_1$ , respectively and 1000 Monte Carlo realizations are used to evaluate the denominator in equation (3). The posterior samples capture well the nonlinear correlation among the parameters and generate a fairly accurate posterior predictive distribution that overlaps with the observations. Additional details can be found in [12, 13].

#### B.5 Friedman 1 model (AdaAnn)

We consider a modified version of the Friedman 1 dataset [22] to examine the performance of our adaptive annealing scheduler in a high-dimensional context. According to the original model in [22], the data are generated as

$$y_i = \mu_i(\boldsymbol{\beta}) + \epsilon_i, \quad \text{where } \mu_i(\boldsymbol{\beta}) = \beta_1 \sin(\pi x_{i,1} x_{i,2}) + \beta_2 (x_{i,3} - \beta_3)^2 + \sum_{j=4}^{10} \beta_j x_{i,j}, \quad (10)$$

where  $\epsilon_i \sim \mathcal{N}(0, 1)$ . We made a slight modification to the model in (10) as

$$\mu_i(\boldsymbol{\beta}) = \beta_1 \sin(\pi x_{i,1} x_{i,2}) + \beta_2^2 (x_{i,3} - \beta_3)^2 + \sum_{j=4}^{10} \beta_j x_{i,j}, \quad (11)$$

and set the true parameter combination to  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{10}) = (10, \pm\sqrt{20}, 0.5, 10, 5, 0, 0, 0, 0, 0)$ . Note that both (10) and (11) contain linear, nonlinear, and interaction terms of the input variables  $X_1$  to  $X_{10}$ , five of which ( $X_6$  to  $X_{10}$ ) are irrelevant to  $Y$ . Each  $X$  is drawn independently from  $\mathcal{U}(0, 1)$ . We used R package `tgp` [23] to generate a Friedman 1 dataset with a sample size of  $n=1000$ . We impose a non-informative uniform prior  $p(\boldsymbol{\beta})$  and, unlike the original modal, we now expect a bimodal posterior distribution of  $\boldsymbol{\beta}$ . Results in terms of marginal statistics and their convergence for the mode with positive  $z_{K,2}$  are illustrated in Table 1 and Figure 5.

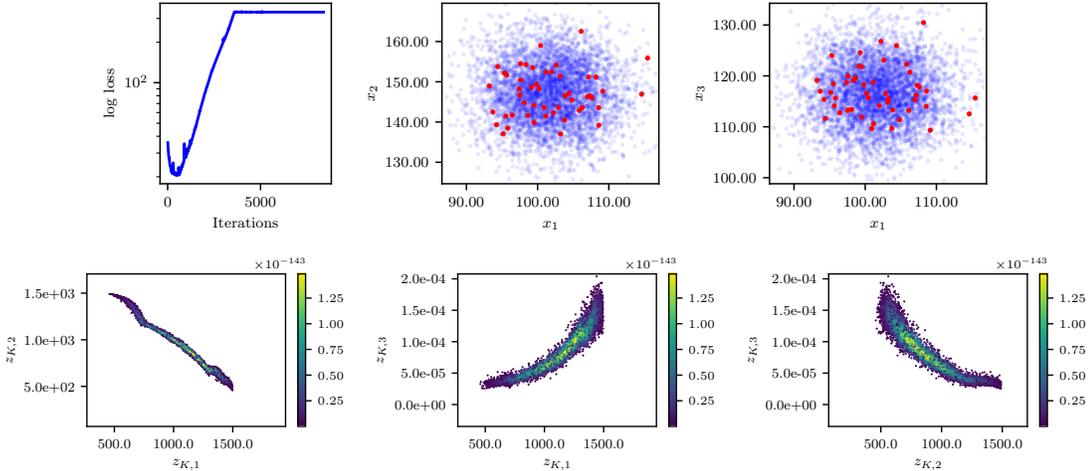


Figure 4: Results from the RCR model. Loss profile (left), posterior predictive distribution (center), and posterior samples (right).

True Value	Mode 1	
	Post. Mean	Post. SD
$\beta_1 = 10$	10.0285	0.1000
$\beta_2 = \pm\sqrt{20}$	4.2187	0.1719
$\beta_3 = 0.5$	0.4854	0.0004
$\beta_4 = 10$	10.0987	0.0491
$\beta_5 = 5$	5.0182	0.1142
$\beta_6 = 0$	0.1113	0.0785
$\beta_7 = 0$	0.0707	0.0043
$\beta_8 = 0$	-0.1315	0.1008
$\beta_9 = 0$	0.0976	0.0387
$\beta_{10} = 0$	0.1192	0.0463

Table 1: Posterior mean and standard deviation for positive mode in the modified Friedman test case.

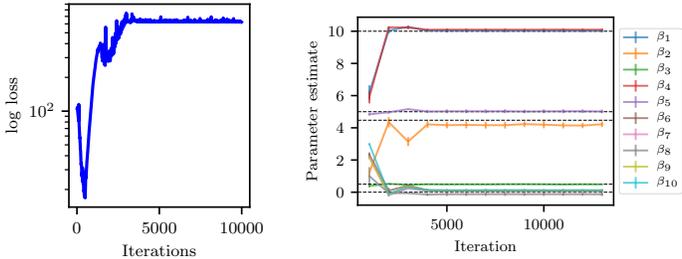


Figure 5: Loss profile (left) and posterior marginal statistics (right) for positive mode in the modified Friedman test case.

### C Hyperparameters in LINFA

This section contains the list of all hyperparameters in the library, their default values, and a description of the functionalities they control. General hyperparameters are listed in Table 6, those related to the optimization process in Table 5, and to the output folder and files in Table 2. Hyperparameters for the proposed NoFAS and AdaAnn approaches are listed in Table 3 and 7, respectively. Finally, a hyperparameter used to select the hardware device is described in Table 4.

Table 2: Output parameters

Option	Type	Description
<i>output_dir</i>	string	Name of output folder where results files are written.
<i>log_file</i>	string	Name of the log file which stores the iteration number, annealing temperature, and value of the loss function at each iteration.
<i>seed</i>	int	Seed for random number generator.

Table 3: Surrogate model parameters (NoFAS)

Option	Type	Description
<i>n_sample</i>	int	Batch size used to generate results after <i>save_interval</i> iterations.
<i>calibrate_interval</i>	int	Number of NF iteration between successive updates of the surrogate model (default <i>1000</i> ).
<i>budget</i>	int	Maximum allowable number of true model evaluations.
<i>surr_pre_it</i>	int	Number of pre-training iterations for surrogate model (default <i>40000</i> ).
<i>surr_upd_it</i>	int	Number of iterations for the surrogate model update (default <i>6000</i> ).
<i>surr_folder</i>	string	Folder where the surrogate model is stored (default <i>./</i> ).
<i>use_new_surr</i>	bool	Start by pre-training a new surrogate and ignore existing surrogates (default <i>True</i> ).
<i>store_surr_interval</i>	int	Save interval for surrogate model ( <i>None</i> for no save, default <i>None</i> ).

Table 4: Device parameters

Option	Type	Description
<i>no_cuda</i>	bool	Do not use GPU acceleration.

Table 5: Optimizer and learning rate parameters

Option	Type	Description
<i>optimizer</i>	string	Type of SGD optimizer (default <i>'Adam'</i> ).
<i>lr</i>	float	Learning rate (default <i>0.003</i> ).
<i>lr_decay</i>	float	Learning rate decay (default <i>0.9999</i> ).
<i>lr_scheduler</i>	string	Type of learning rate scheduler ( <i>'StepLR'</i> or <i>'ExponentialLR'</i> ).
<i>lr_step</i>	int	Number of steps before learning rate reduction for the step scheduler.
<i>log_interval</i>	int	Number of iterations between successive loss printouts (default <i>10</i> ).

Table 6: General parameters

Option	Type	Description
<i>name</i>	str	Name of the experiment.
<i>flow_type</i>	str	type of normalizing flow ( <i>'maf'</i> , <i>'realnvp'</i> ).
<i>n_blocks</i>	int	Number of normalizing flow layers (default 5).
<i>hidden_size</i>	int	Number of neurons in MADE and RealNVP hidden layers (default 100).
<i>n_hidden</i>	int	Number of hidden layers in MADE (default 1).
<i>activation_fn</i>	str	Activation function for MADE network used by MAF (default <i>'relu'</i> ).
<i>input_order</i>	str	Input order for MADE mask creation ( <i>'sequential'</i> or <i>'random'</i> , default <i>'sequential'</i> ).
<i>batch_norm_order</i>	bool	Adds batchnorm layer after each MAF or RealNVP layer (default <i>True</i> ).
<i>save_interval</i>	int	How often to save results from the normalizing flow iterations. Saved results include posterior samples, loss profile, samples from the posterior predictive distribution, observations, and marginal statistics.
<i>input_size</i>	int	Input dimensionality (default 2).
<i>batch_size</i>	int	Number of samples from the basic distribution generated at each iteration (default 100).
<i>true_data_num</i>	int	Number of additional true model evaluations at each surrogate model update (default 2).
<i>n_iter</i>	int	Total number of NF iterations (default 25001).

Table 7: Parameters for the adaptive annealing scheduler (AdaAnn)

Option	Type	Description
<i>annealing</i>	bool	Flag to activate the annealing scheduler. If this is <i>False</i> , the target posterior distribution is left unchanged during the iterations.
<i>scheduler</i>	string	Type of annealing scheduler ( <i>'AdaAnn'</i> or <i>'fixed'</i> , default <i>'AdaAnn'</i> ).
<i>tol</i>	float	KL tolerance. It is kept constant during inference and used in the numerator of equation (3).
<i>t0</i>	float	Initial inverse temperature.
<i>N</i>	int	Number of batch samples during annealing.
<i>N_1</i>	int	Number of batch samples at $t = 1$ .
<i>T_0</i>	int	Number of initial parameter updates at $t_0$ .
<i>T</i>	int	Number of parameter updates after each temperature update. During such updates the temperature is kept fixed.
<i>T_1</i>	int	Number of parameter updates at $t = 1$
<i>M</i>	int	Number of Monte Carlo samples used to evaluate the denominator in equation (3).

## References

- [1] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092, 1953.
- [2] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications. 1970.
- [3] Stuart Geman and Donald Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence*, (6):721–741, 1984.

- [4] Alan E Gelfand and Adrian FM Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association*, 85(410):398–409, 1990.
- [5] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [6] Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [7] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- [8] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. *arXiv preprint arXiv:1605.08803*, 2016.
- [9] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. *Advances in neural information processing systems*, 29: 4743–4751, 2016.
- [10] Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- [11] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(1):2617–2680, 2021.
- [12] Yu Wang, Fang Liu, and Daniele E Schiavazzi. Variational inference with nofas: Normalizing flow with adaptive surrogate for computationally expensive models. *Journal of Computational Physics*, 467:111454, 2022.
- [13] Emma R Cobian, Jonathan D Hauenstein, Fang Liu, and Daniele E Schiavazzi. Adaann: Adaptive annealing scheduler for probability density approximation. *International Journal for Uncertainty Quantification*, 13, 2023.
- [14] Bingyue Su, Yu Wang, Daniele E Schiavazzi, and Fang Liu. Differentially private normalizing flows for density estimation, data synthesis, and variational inference with application to electronic health records. *arXiv preprint arXiv:2302.05787*, 2023.
- [15] Ali Siahkoohi, Gabrio Rizzuti, Mathias Louboutin, Philipp Witte, and Felix Herrmann. Preconditioned training of normalizing flows for variational inference in inverse problems. In *Third Symposium on Advances in Approximate Bayesian Inference*, 2021. URL <https://openreview.net/forum?id=P9m1sMaNq8T>.
- [16] Michael Brennan, Daniele Bigoni, Olivier Zahm, Alessio Spantini, and Youssef Marzouk. Greedy inference with structure-exploiting lazy maps. *Advances in Neural Information Processing Systems*, 33:8330–8342, 2020.
- [17] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- [18] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation, 2018.
- [19] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning*, pages 881–889. PMLR, 2015.
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [21] Ilya M Sobol’. Theorems and examples on high dimensional model representation. *Reliability Engineering and System Safety*, 79(2):187–193, 2003.
- [22] Jerome H Friedman. Multivariate adaptive regression splines. *The annals of statistics*, 19(1):1–67, 1991.
- [23] Robert B Gramacy. tgp: an r package for bayesian nonstationary, semiparametric nonlinear regression and design by treed gaussian process models. *Journal of Statistical Software*, 19:1–46, 2007.