

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

Bayesian neural networks via MCMC: a Python-based tutorial

ROHITASH CHANDRA^{1,2}, (SM, IEEE), and Joshua Simmons³

¹Transitional Artificial Intelligence Research Group, School of Mathematics and Statistics, UNSW Sydney, Australia (e-mail: rohitash.chandra@unsw.edu.au)

²Pingala Institute of Artificial Intelligence, Sydney, Australia

³ARC Training Centre in Data Analytics for Resources and Environments (DARE), University of Sydney, Australia (e-mail: joshua.simmons@sydney.edu.au)

Corresponding author: First A. Author (e-mail: rohitash.chandra@unsw.edu.au).

This work was supported in part by the Australian Research Council through Grant IC190100031. We thank Royce Chen from UNSW Sydney for his contribution to the earlier version of this paper.

ABSTRACT Bayesian inference provides a methodology for parameter estimation and uncertainty quantification in machine learning and deep learning methods. Variational inference and Markov Chain Monte-Carlo (MCMC) sampling methods are used to implement Bayesian inference. In the past three decades, MCMC sampling methods have faced some challenges in being adapted to larger models (such as in deep learning) and big data problems. Advanced proposal distributions that incorporate gradients, such as a Langevin proposal distribution, provide a means to address some of the limitations of MCMC sampling for Bayesian neural networks. Furthermore, MCMC methods have typically been constrained to statisticians and currently not well-known among deep learning researchers. We present a tutorial for MCMC methods that covers simple Bayesian linear and logistic models, and Bayesian neural networks. The aim of this tutorial is to bridge the gap between theory and implementation via coding, given a general sparsity of libraries and tutorials to this end. This tutorial provides code in Python with data and instructions that enable their use and extension. We provide results for some benchmark problems showing the strengths and weaknesses of implementing the respective Bayesian models via MCMC. We highlight the challenges in sampling multi-modal posterior distributions for the case of Bayesian neural networks and the need for further improvement of convergence diagnosis methods.

INDEX TERMS MCMC; Bayesian deep learning; Bayesian neural networks; Bayesian linear regression; Bayesian inference

A. INTRODUCTION

BAYESIAN inference provides a probabilistic approach for parameter estimation in a wide range of models used across the fields of machine learning, econometrics, environmental and Earth sciences [1]–[5]. The term ‘probabilistic’ refers to the representation of unknown parameters as probability distributions rather than using fixed point estimates as in conventional machine learning models where gradient-based optimisation methods are prominent [6]. A probabilistic representation of unknown parameters requires a different approach to optimisation, which is known as *sampling* from a computational statistics point-of-view [7].

Markov Chain Monte-Carlo (MCMC) sampling methods have been prominent for inference (estimation) of model parameters via the posterior probability distribution. In other words, Bayesian methods attempt to quantify the uncertainty in model parameters by marginalising over the predictive

posterior distribution. Hence, in the case of neural networks, MCMC methods can be used to implement Bayesian neural networks that represent weights and biases as probability distributions [8]–[12]. Probabilistic machine learning provides natural way of providing uncertainty quantification in predictions [13], since the uncertainties can be obtained by probabilistic representation of parameters. This inference procedure can be seen as a form of learning (optimisation) applied to the model parameters [9]. In this tutorial, we employ linear models and simple neural networks to demonstrate the use of MCMC sampling methods. The probabilistic representation of weights and biases in the respective models allows uncertainty quantification on model predictions.

We note that MCMC refers to a family of algorithms for implementing Bayesian inference for parameter and uncertainty estimation in models. Bayesian inference applications include statistical, graphical, and machine learning models. The differences in the model complexity from dif-

ferent domains have led to the existence of a wide range of MCMC sampling algorithms. Some of the prominent ones are Metropolis-Hastings algorithm [14]–[16], Gibbs sampling [17]–[19], Langevin MCMC [20]–[22], rejection sampling [23], [24], importance sampling [25], [26], sequential MCMC [27], adaptive MCMC [28], parallel tempering (tempered) MCMC [29]–[32], reversible-jump MCMC [33], [34], specialised MCMC methods for discrete time series models [35]–[37], constrained parameter and model settings [38], [39], and likelihood free MCMC [40]. MCMC sampling methods have also been used for data augmentation [41], [42], model fusion [43], model selection [44], [45], and interpolation [46]. Apart from this, we note that MCMC methods have been prominent in a wide range of applications that include geophysical inversions [47]–[49], geoscientific models [5], [50], [51], environmental and hydrological modelling [52], [53], bio-systems modelling [54]–[56], and quantitative genetics [57], [58].

In the case of Bayesian neural networks, the large number of model parameters that emerge from large neural network architectures and deep learning models pose challenges for MCMC sampling methods. Hence, progress in the application of Bayesian approaches to big data and deep neural networks has been slow. Research in this space has included a number of methods that have been fused with MCMC such as gradient-based methods [22], [59]–[62], and evolutionary (meta-heuristic) algorithms which include differential evolution, genetic algorithms, and particle swarm optimisation [63]–[66].

This use of gradients in MCMC was initially known as Metropolis-adjusted Langevin dynamics [22] and has shown promising performance for linear models [61] and has also been extended to Bayesian neural networks [62]. Hamiltonian Monte Carlo (HMC) sampling also employ gradient-based proposal distributions [60] and has been effectively applied to Bayesian neural networks [67]. In similar way, Langevin dynamics can be used to incorporate gradient-based stepping with Gaussian noise into the proposal distribution [61]. HMC avoids random walk behaviour using an auxiliary momentum vector and implementing Hamiltonian dynamics where the momentum samples are discarded later. The samples are hence less correlated and tend converge to the target distribution more rapidly. Another direction has been the use of better exploration features in MCMC sampling such as parallel tempering MCMC with Langevin proposal distribution and parallel computing [62]. These have the ability to provide a competitive alternative to stochastic gradient-descent [68] and Adam optimizers [6] with the addition of uncertainty quantification in predictions. These methods have also been applied to Bayesian deep learning models such as Bayesian autoencoders [69] and Bayesian graph convolutional neural networks (CNNs) [70] which require millions of trainable parameters to be represent as posterior distributions. Recently, Kapoor et al. [66] combined tempered MCMC with particle swarm optimisation-based proposal distribution in a parallelized environment that

showed more effective sampling when compared with the conventional approach. However, we note that large deep learning models can feature hundreds of millions to billions of parameters, which brings further challenges to sampling strategies and hence the road is less travelled.

Variational inference provides an alternative approach to MCMC methods to approximate Bayesian posterior distribution [71], [72]. *Bayes by backpropagation* is a VI method that showed competitive results when compared to stochastic gradient descent and *dropout* methods used as approximate Bayesian methods [73]. Dropout is a regularisation technique that involves randomly dropping selected weights in forward-pass operation of backpropagation. This improves the generalization performance of neural networks and has been widely adopted [74]. Gal and Ghahramani [75] presented an approximate Bayesian methodology based on dropout-based regularisation which has been used for other deep learning models such as CNNs [76]. Later, Gal and Ghahramani [77] presented variational inference-based dropout technique for recurrent neural networks (RNNs); particularly, long-short term memory (LSTM) and gated recurrent unit (GRU) models for language modelling and sentiment analysis tasks.

We argue that the use of dropouts for Bayesian inference [76] cannot be seen as an alternative to MCMC sampling which samples directly from the posterior distribution. In the case of dropouts for Bayesian inference, we do not know the priors nor know much about the posterior distribution and there is little theoretical rigour, only computational efficacy or capturing noise and uncertainty during model training. Furthermore, in the Bayesian methodology, a probabilistic representation using priors is needed which is questionable in the dropout methodology for Bayesian computation. Given that variational inference methods are seen as approximate Bayesian methods, we need to invest more effort in directly sampling from the posterior distribution for Bayesian deep learning models. This can only be possible if both communities (i.e., statistics and machine learning) are aware about the strengths and weaknesses of MCMC methods for sampling Bayesian neural networks that span hundreds to thousands of parameters, and go orders of magnitude higher when looking at Bayesian deep learning models. The progress of MCMC for deep learning has been slow, due to lack of implementation details, libraries and tutorials that provide that balance of theory and implementation.

In this paper, we present a Python-based MCMC sampling tutorial for simple Bayesian linear models and Bayesian neural networks. We provide code in Python with data and instructions that enable their use and extension. We provide detailed instructions for sample code in a related Github repository which is easy to clone and run. Our code implementation is simple and relies on basic Python libraries such as *numpy* as the goal of this tutorial is to serve as a go-to document for beginners who have basic knowledge of machine learning models, and need to get hands-on experience with MCMC sampling. Hence, this is a code-based computational tutorial with a theoretical background. We provide results for some

benchmark problems showing the strengths and weaknesses of implementing the respective Bayesian models. Finally, we highlight the challenges in sampling multi-modal posterior distributions in the case of Bayesian neural networks and shed light on the use of convergence diagnostics.

The rest of the paper is organised as follows. In Section 2, we present a background and literature review of related methods. Section 3 presents the proposed methodology that includes MCMC sampling code implementation in Python for respective models, followed by experiments and results in Section 4. Section 5 provides a discussion and Section 6 concludes the paper with directions for future work.

B. BACKGROUND

1) Bayesian inference

We recall that Bayesian methods account for the uncertainty in prediction and decision-making via the posterior distribution [78]. Note that the posterior is the conditional probability determined after taking into account the prior distribution and the relevant evidence or data via sampling methods. Thomas Bayes (1702 – 1761) presented and proved a special case of the Bayes' theorem [79], [80] which is the foundation of Bayesian inference. However, it was Pierre-Simon Laplace (1749 – 1827) who introduced a general version of the theorem and used it to approach problems [81]. Figure 1 gives an overview of the Bayesian inference framework that uses data with a prior and likelihood to construct to sample from the posterior distribution. This is the building block of the rest of the lessons that will feature Bayesian logistic regression and Bayesian neural networks.

Bayesian inference estimates unknown parameters using prior information or belief about the variable. Prior information is captured in the form of a distribution. A simple example of a prior belief is a distribution that has a positive real-valued number in some range. This essentially would imply a belief that our result or posterior distribution would likely be a distribution of positive numbers in some range which would be similar to the prior but not the same. If the posterior and prior both follow the same type of distribution, this is known as a *conjugate prior* [82]. If the prior provides useful information about the variable, rather than very loose constraints, it is known as an *informative prior*. The prior distribution is based on expert knowledge (opinion) and is also dependent on the domain for different types of models [83], [84].

The need for efficient sampling methods to implement Bayesian inference has been a significant focus of research in computational statistics. This is especially true in the case of multimodal and irregular posterior distributions [32], [85], [86] which tend to dominate Bayesian neural networks [11], [87]. MCMC sampling methods are used to update the probability for a hypothesis (proposal Θ) as more information becomes available. The hypothesis is given by a prior probability distribution that expresses one's belief about a quantity (or free parameter in a model) before some data (d) are observed. MCMC sampling methods construct the

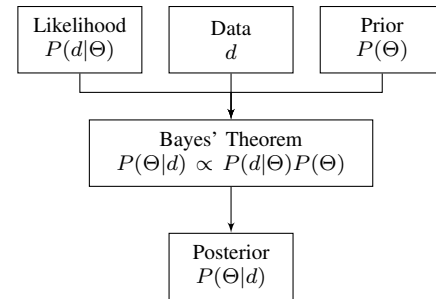


FIGURE 1: We show the relationship of the likelihood with data and prior distribution for sampling the posterior distribution.

posterior distribution ($P(\Theta|d)$) iteratively, using a proposal distribution, prior distribution $P(\Theta)$ and a likelihood function ($P(d|\Theta)$), as expressed below

$$P(\Theta|d) = \frac{P(d|\Theta) \times P(\Theta)}{P(d)}. \quad (1)$$

$P(d)$ is the marginal distribution of the data and is often seen as a normalising constant and ignored. Hence, ignoring it, we can also express the above by

$$P(\Theta|d) \propto P(d|\Theta) \times P(\Theta). \quad (2)$$

The likelihood function is a function of the parameters of a given model provided specific observed data [88]. The likelihood function can be seen as a measure of fit to the data, given the proposals that are drawn from the proposal distribution. Hence, from an optimisation perspective, the likelihood function can be seen as a fitness or error function. The posterior distribution is constructed after taking into account the relevant evidence (data) and prior distribution, with the likelihood that considers the proposal and the model. MCMC methods essentially implement Bayesian inference via a numerical approach that marginalizes or integrates over the posterior distribution [89]. Note that *probability* and *likelihood* are not the same in the field of statistics, while in everyday language they are used as if they are the same. The term "probability" refers to the possibility of something happening, in relation to a given distribution of data. The likelihood refers to the likelihood function that provides a measure of fitness in relation to a distribution. The likelihood function indicates which parameter (data) values are more likely than others in relation to a distribution. Further and detailed explanations regarding Bayesian inference and MCMC sampling have been given in the literature [90], [91].

2) Probability distributions

a: Gaussian (Normal) Distribution

A normal probability density or distribution, also known as the Gaussian distribution, is described by two parameters, mean (μ) around which the distribution is centered and the standard deviation (σ) which describes the spread (sometimes described instead by the variance, σ^2). Using these two

parameters, we can fit a probability (normal) distribution to data from some source. In a similar way, given a probability distribution, we can generate data and this process is known as sampling from the distribution. In sampling the distribution, we simply present random data points (uniform) to the distribution and get data that are, in a way, transformed by the distribution. These parameters determine the shape of the probability distribution, e.g., if it is peaked or spread. Note that the normal distribution is symmetrical and caters for negative and positive numbers of real data.

Equation 3 presents the Gaussian distribution probability density function (PDF) for parameters μ and σ .

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2\right) \quad (3)$$

We will sample from this distribution in Python via the NumPy library [92] which covers the various distributions discussed in this tutorial. The SciPy library [93] is used to get a representation of the probability distribution function (PDF). The associated Github repository¹ contains the code to generate Figures 2 to 4 using the Seaborn and Matplotlib Python libraries. Listing 1 shows an example of drawing samples from a Gaussian distribution and obtaining the PDF values across the domain 0 to 8. The SciPy and NumPy libraries cover all distributions mentioned below.

```
1 import numpy as np
2 from numpy import random
3 from scipy import stats
4
5 # 10 random draws from a standard Gaussian
6 # distribution - N(0,1)
7 samples = random.randn(10)
8 # Gaussian distribution PDF with mean of 4 and
9 # standard deviation of 0.5
10 x = np.linspace(0, 8, 100)
11 pdf = stats.norm.pdf(x, loc=4, scale=0.5)
```

Listing 1: Random number generation for a Gaussian distribution

We note that the mean and standard deviation are purely based on the data that will change depending on the dataset. Let us visualise what happens to the distribution when the mean and standard deviation change, as shown in Figure 2.

b: Multivariate Normal distribution

The multivariate normal distribution or joint normal distribution generalises univariate normal distribution to more variables or higher dimensions, as shown in the PDF in Equation 4.

$$f(x_1, \dots, x_M) = \frac{1}{\sqrt{(2\pi)^M |\Sigma|}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right) \quad (4)$$

¹<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/01-Distributions.ipynb>

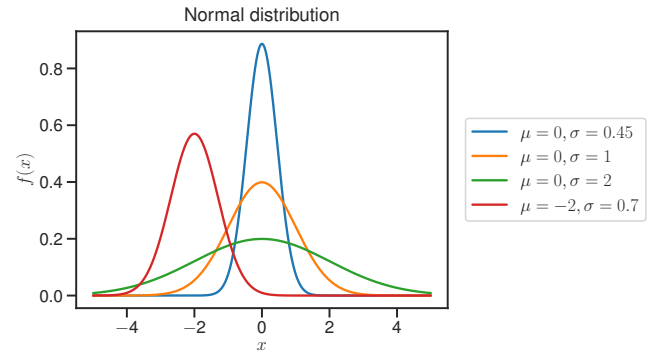


FIGURE 2: Normal distributions with different parameters, i.e., mean and the standard deviation.

where \mathbf{x} is a real M -dimensional column vector and $|\Sigma|$ is the determinant of the *symmetric covariance matrix*, which is *positive definite*.

c: Gamma distribution

A gamma distribution is defined by the parameters shape (α) and rate (β), as shown below.

$$f(x; \alpha, \beta) = \frac{\beta^\alpha x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)} \quad (5)$$

for $x > 0$, $\alpha, \beta > 0$; where $\Gamma(n) = (n-1)!$. Figure 3 presents the Gamma distribution for various parameter combinations, with the corresponding code in the accompanying Github repository.

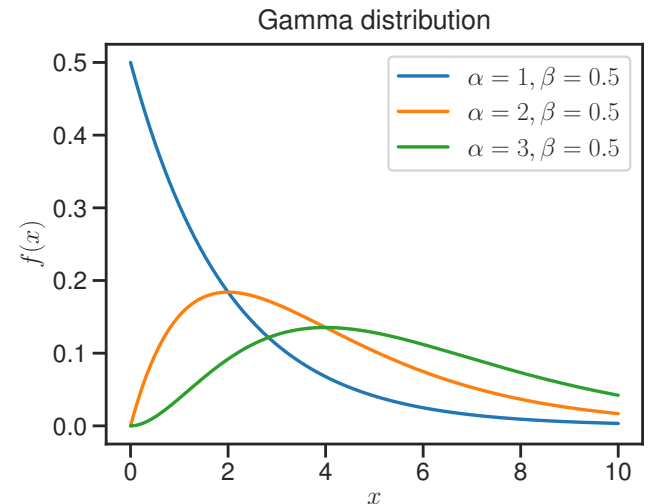


FIGURE 3: Gamma distributions with different shape and rate parameters (α and β).

The corresponding inverse-Gamma (IG) distribution takes the same parameters with examples given in Figure 4 and is more appropriate for real positive numbers.

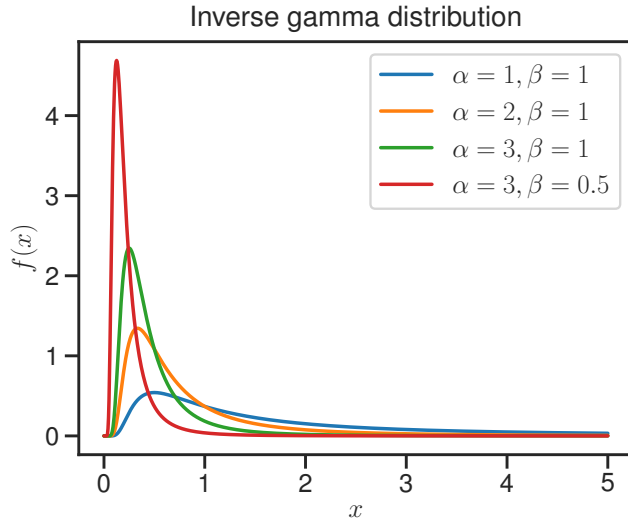


FIGURE 4: Inverse gamma distributions with different shape and rate parameters (α and β).

d: Binomial distribution

So far, we have only addressed real numbers with respective probability distributions; however, we also need to consider discrete numbers. The Bernoulli distribution is a discrete probability distribution typically used for modelling binary classification problems. We begin with an example where a variable x takes the value 1 with probability p and the value 0 with probability $q = 1 - p$. We give the probability mass function for this distribution over the possible outcomes (x) in Equation 6.

$$f(x; p) = p^x (1 - p)^{1-x} \quad (6)$$

for $x \in \{0, 1\}$. The probability of getting exactly k successes ($x = 1$) in n independent Bernoulli trials ($f(k, n, p)$) is given as $\Pr(k; n, p)$ in Equation 7.

$$\Pr(k; n, p) = \binom{n}{k} p^k (1 - p)^{n-k} \quad (7)$$

for $k = 0, 1, 2, \dots, n$, where $\binom{n}{k} = \frac{n!}{k!(n-k)!}$.

e: Multinomial distribution

Previously, we catered for the case of two outcomes; however, we can consider the case of more than two outcomes. Suppose a single trial can result in k ($k \geq 2$) possible outcomes numbered $1, 2, \dots, k$ and let $p_i = \mathbb{P}(\text{a single trial results in outcome } i)$ ($\sum_{i=1}^k p_i = 1$). In the case of n independent trials, let X_i denote the number of trials resulting in outcome i (then $\sum_{i=1}^k X_i = n$). Then, we can state that the distribution of $(X_1, X_2, \dots, X_k) \sim \text{Multinomial}(n; p_1, p_2, \dots, p_k)$, and it holds

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}, \quad 0 < p_i < 1, \sum_{i=1}^k p_i = 1. \quad (8)$$

C. MCMC

We begin by noting that a Markov process is uniquely defined by its transition probabilities $P(x'|x)$, which defines the probability of transitioning from any given state x to another given state x' . The Markov process has a unique stationary distribution $\pi(x)$ given the following two conditions are met.

- 1) There must exist a stationary distribution π which solves the detailed balance equations, and therefore requires that each transition $x \rightarrow x'$ is reversible. This implies that for every pair of states x, x' , the probability of being in state x and moving to state x' , must be equal to the probability of being in state x' and moving to state x ; hence, $\pi(x)P(x' | x) = \pi(x')P(x | x')$.
- 2) The stationary distribution must be unique, which is guaranteed by ergodicity of the Markov process [94]–[97]. Ergodicity is guaranteed when every state is aperiodic (i.e., the system does not return to the same state at fixed intervals) and positive recurrent (i.e., the expected number of steps for returning to the same state is finite). An ergodic system is one that mixes well, in other words, you get the same result whether you average its values over time or space.

Given that $\pi(x)$ is chosen to be $P(x)$, the condition of detailed balance becomes $P(x' | x)P(x) = P(x | x')P(x')$ which is re-written as shown in Equation 9.

$$\frac{P(x' | x)}{P(x | x')} = \frac{P(x')}{P(x)} \quad (9)$$

Algorithm 1 presents a basic MCMC sampler with random-walk proposal distribution that runs until a maximum number of samples (N_{max}) has been reached for training data, \mathbf{d} .

Algorithm 1 proceeds by proposing new values of the parameter x (Step 1) from the selected proposal distribution $q(\cdot)$; in this case, a uniform distribution between 0 and 1. Conditional on these proposed values, the model $f(x', \mathbf{d})$ computes or predicts an output using proposal x' and data \mathbf{d} (Step 2). We compute the likelihood using the prediction and employ a Metropolis-Hasting criterion (Step 3) to determine whether to accept or reject the proposal (Step 5). We compare the acceptance ratio α with $u \sim U(0, 1)$, this enforces that the proposal is accepted with probability α . If the proposal is accepted, the chain moves to this proposed value. If rejected, the chain stays at the current value. The process is repeated until the convergence criterion is met, which is the maximum number of samples (N_{max}).

1) Priors

The prior distribution is generally based on belief, expert opinion or other information without viewing the data [9],

Data: Training data, \mathbf{d}

Result: N_{max} samples from the posterior distribution

- Initialise x_0 ;

for $i = 1$ **until** N_{max} **do**

1. Propose a value $x'|x_i \sim q(x_i)$, where $q(\cdot)$ is the proposal distribution;

2. Given x' , execute the model $f(x', \mathbf{d})$ to compute the predictions (output y) and the likelihood;

3. Calculate the acceptance probability
 $\alpha = \min \left(1, \frac{P(x') q(x_i|x')}{P(x_i) q(x'|x_i)} \right)$

4. Generate a random value from a uniform distribution $u \sim U(0, 1)$;

5. Accept or reject proposed value x' ;

if $u < \alpha$ **then**

| accept the sample, $x_i = x'$

else

| reject current and retain previous sample,
 $x_i = x_{i-1}$

end

end

Algorithm 1: A basic MCMC sampler leveraging the Metropolis-Hastings algorithm

[98]. Information to construct the prior can be based on past experiments or the posterior distribution of the model for related datasets. There are no hard rules for how much information should be encoded in the prior distribution; hence, we can take multiple approaches.

An *informative prior* gives specific and definite information about a variable. If we consider the prior distribution for the temperature tomorrow evening, it would be reasonable to use a normal distribution with an expected value (as mean) of today's evenings temperature with a standard deviation of the temperature each evening for the entire season. A *weakly informative prior* expresses partial information about a variable. In the case of the prior distribution of evening temperature, a weakly informative prior would consider day time temperature of the day (as mean) with a standard deviation of day time temperature for the whole year. An *uninformative prior* or *diffuse prior* expresses vague information about a variable, such as the variable is positive or has some limit range.

A number of studies have been done regarding priors for linear models [99], [100] and Bayesian neural networks and deep learning models [101]. Hobbs et al. [102] presented a study for Bayesian priors in generalised linear models for clinical trials. We note that incorporation of prior knowledge in deep learning models [103], is different from selecting or defining priors in Bayesian deep learning models. Due to the similarity of terms, we caution the readers that these can be often confused and mixed up.

In the case of Bayesian neural networks, the prior distribution can be based on the distribution of the weights and biases from similar neural network models. This can be seen as an

example of expert knowledge and implemented in previous studies [69], [104]. Another example of expert knowledge is the concept of *weight decay* [105] regularisation (L2 or Ridge regression [106]) which restricts large weights and can be incorporated when defining the prior distribution (priors) [8], [8], [9].

2) MCMC sampler in Python

We begin with a deliberately simple example where we sample one parameter from a binomial distribution to demonstrate a simple MCMC implementation in Python. Looking at a simple binomial (e.g., coin flipping) likelihood (we will explore the likelihood later), given the data of k successes in n trials, we calculate the posterior probability of the parameter p that defines the chance of success for any given trial. MCMC sampling requires a prior distribution and a likelihood function to evaluate a set of parameters (proposed) for the given data and model. In other words, the likelihood is the measure of the quality of proposals obtained from a proposal distribution.

Listing 2 presents an implementation ² of this simple MCMC sampling exercise in Python of Algorithm 1.

In this example, we adopt a uniform distribution as an uninformative prior, only constraining the p to be between the values of 0 and 1 ($p \in [0, 1]$).

In MCMC sampling, a certain portion of the initial samples are discarded that are known as the burn-in or warmup period which could be seen as an optimisation stage. The burn-in period depends on the sampling problem (complexity of the model). In this simple case, we will use 25 % burn-in and in the case of neural network models, 50 % burn-in will likely be required. Essentially, during burn-in, we are discarding material that is not part of the posterior distribution, since the posterior should feature good proposals which we get once the sampler goes towards convergence.

Typically, histograms of the posterior distribution and the trace plot are used to visualise the MCMC sampling performance. The histogram of the posterior distribution allows us to examine the mean and variance visually, while the trace plot shows the value of samples at each iteration, allowing us to examine the behaviour and convergence of the MCMC.

Although it is necessary to exclude the burn-in samples in the posterior distribution, it can be helpful to include them in the trace plots so that we can examine where the model started and how well it converged. We provide a visualisation of the results that features a normal distribution obtained by a simple MCMC sampling from executing code Listing 2. The histogram of the posterior shows a normally distributed shape (Figure 5 - Panel a), and the trace plot (Figure 5 - Panel b) shows that the samples are distributed around the convergence value, as well as the burn-in samples which are in red. We also note that the value of the posterior is usually taken as the mean of the distribution, and in this

²<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/02-Basic-MCMC.ipynb>

```

1 # First define our likelihood function which will be dependent on provided 'data'
2 # in this case we will choose k = 50, n = 100
3 def likelihood(query_prob):
4     """
5     Given the data of k successes in n trials, return a likelihood function which
6     evaluates the probability that a single success (p) is query_prob for any given
7     query_prob (between 0 and 1).
8     """
9     k = 50
10    n = 100
11    return stats.binom.pmf(k, n, query_prob)
12
13 ## MCMC Settings and Setup
14 n_samples = 10000 # number of samples to draw from the posterior
15 burn_in = 2500 # number of samples to discard before recording draws from the posterior
16
17 x = random.uniform(0, 1) # initialise a value of x0
18
19 # create an array of NaNs to fill with our samples
20 p_posterior = np.full(n_samples, np.nan)
21
22 print('Generating {} MCMC samples from the posterior:'.format(n_samples))
23
24 # now we can start the MCMC sampling loop
25 for ii in np.arange(n_samples):
26     # Sample a value uniformly from 0 to 1 as a proposal
27     x_new = random.uniform(0, 1)
28
29     # Calculate the Metropolis-Hastings acceptance probability based on the prior
30     # (can be ignored in this case) and likelihood
31     prior_ratio = 1 # for this simple example as discussed above
32     likelihood_ratio = likelihood_function(x_new) / likelihood_function(x)
33     alpha = np.min([1, likelihood_ratio * prior_ratio])
34
35     # Here we use a random draw from a uniform distribution between 0 and 1 as a
36     # method of accepting the new proposal with a probability of alpha
37     # (i.e., accept if u < alpha)
38     u = random.uniform(0, 1)
39     if u < alpha:
40         x = x_new # then update the current sample to the propoal for the next iteration
41
42     # Store the current sample
43     p_posterior[ii] = x

```

Listing 2: Python implementation of Algorithm 1

case, the mean value is 0.502. The median can also be taken as a measure which would be more useful in irregular distributions.

D. BAYESIAN LINEAR MODELS VIA MCMC

We provide details of implementing *Bayesian linear regression* that uses MCMC sampling with random-walk proposal distribution. We wish to model a dataset consisting of input (features or covariates) $\mathbf{x} = (x_1, \dots, x_S)'$ and corresponding outputs $\mathbf{y} = (y_1, \dots, y_S)'$ for S instances in data. This approach models the response observations as being composed of a regression component (the linear regression denoted by $f(\mathbf{x}, \theta)$) and a noise term (Gaussian distribution with a mean of zero ($\mu = 0$) and variance τ^2 (see Equation 10). In the Bayesian linear regression, we treat the parameters (θ and τ^2) as random variables to be estimated (sampled) based on the data and likelihood. So

$$\bar{y} = f(\mathbf{x}, \theta) + e \quad e \sim \mathcal{N}(0, \tau^2) \quad (10)$$

or equivalently,

$$p(\mathbf{y}|\mathbf{x}, \theta, \tau^2) \sim \mathcal{N}(f(\mathbf{x}, \theta), \tau^2). \quad (11)$$

Equation 12 expresses the general case of a linear model using a vector of input data \mathbf{x} to obtain prediction \mathbf{y} .

$$f(\mathbf{x}, \theta) = \theta \mathbf{x}^T \quad (12)$$

In the case of Bayesian linear regression, θ is a set of distributions (typically Gaussian) rather than a fixed point estimate in conventional linear models. Therefore, we estimate the parameters (θ and τ^2) using MCMC sampling to obtain their posterior distributions. Note that from an optimisation perspective, the case of sampling can be seen as a form of optimisation, e.g., using gradient-based methods for learning the parameters of linear models and neural networks in the machine learning and neural networks literature [107], [108]. As noted earlier, the key feature of a MCMC sampler is the ability to sample a posterior probability distribution that

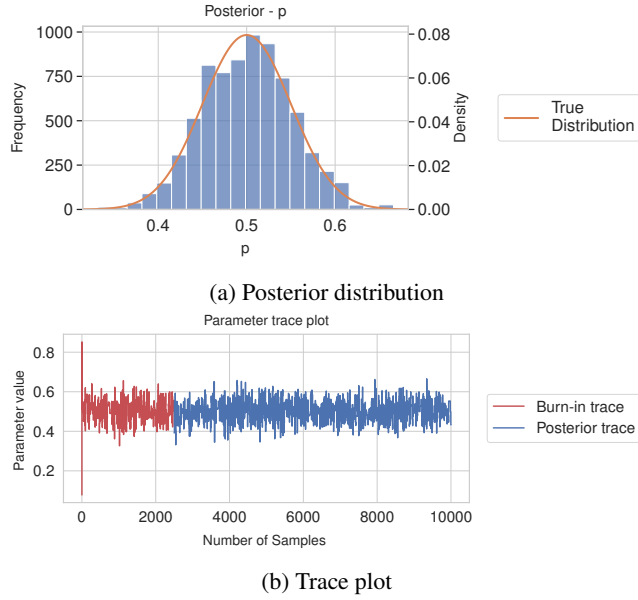


FIGURE 5: Posterior and trace plot for the basic MCMC sampler given in Listing 2.

represents the parameters of a model rather than a fixed point estimate given by optimisation methods.

1) Likelihood

Our Bayesian approach for the problem requires sampling (estimating) the posterior distribution $p(\theta | \mathbf{y})$ that requires the definition of both a likelihood function $p(\theta | \mathbf{x})$ and prior distribution $p(\theta)$. We begin by defining the likelihood function, i.e. probability of the data given the model, which is given by the product of the likelihood for every data point in the dataset of S instances, as shown in Equation 13

$$p(\mathbf{y} | \mathbf{x}, \theta, \tau^2) = \prod_{t=1}^S p(y_t | \mathbf{x}_t, \theta, \tau^2). \quad (13)$$

We note that for our MCMC sampler, we use the log-likelihood (i.e., taking the log of the likelihood function) to eliminate numerical instabilities, which can occur since we multiply multiple probabilities together which grows with the size of the data. It is also more convenient to maximize the log of the likelihood function since the logarithm is monotonically increasing function of its argument, i.e., maximization of the log of a function is equivalent to maximization of the function itself. In order to transform a likelihood function into a log-likelihood, we will use the log product rule as given below

$$\log_b(x \times y) = \log_b(x) + \log_b(y). \quad (14)$$

The log-likelihood simplifies the subsequent mathematical analysis and also helps avoid numerical instabilities due to the product of a large number of small probabilities. In the log-likelihood, Equation 13 is much simplified by computing the sum of the log probabilities as given in Equation 15

$$\ln p(\mathbf{y} | \mathbf{x}, \theta, \tau^2) = \sum_{t=1}^S \ln p(y_t | \mathbf{x}_t, \theta, \tau^2). \quad (15)$$

In order to construct the likelihood function, we use our definition of the probability for each data point given the model as shown in Equation 11, and the form of the Gaussian distribution as defined in Equation 3. We use a set of weights and biases as the model parameters θ in our model $f(\mathbf{x}, \theta)$, for S training data instances and variance τ^2 . Our assumption of normally distributed errors leads to a likelihood given in Equation 16

$$p(\mathbf{y} | \mathbf{x}, \theta, \tau^2) = \frac{1}{(2\pi\tau^2)^{S/2}} \exp \left(-\frac{1}{2\tau^2} \sum_{t=1}^S (y_t - f(\mathbf{x}_t, \theta))^2 \right). \quad (16)$$

2) Prior

We note that a conventional linear model transforms into a Bayesian linear model with the use of a prior distribution and a likelihood function to sample the posterior distribution via the MCMC sampler. In Section -C2, we discussed the need to define a prior distribution for our model parameters θ and τ^2 . In the case where the prior distribution comes from the same family as the posterior distribution, the prior and posterior are then called conjugate distributions [109], [110]. The prior is called a *conjugate prior* for the likelihood function of the Bayesian model.

To implement conjugate priors in our linear model, we will assume a multivariate Gaussian prior for θ (Equation 17) and an inverse Gamma distribution (IG) for τ^2 (Equation 18).

$$\theta \sim \mathcal{N}(0, \sigma^2) \quad (17)$$

$$\tau^2 \sim IG(\nu_1, \nu_2) \quad (18)$$

The noise is defined by the spread (variance) of a normal distribution. We need to define the variance that is represented by τ^2 which cannot be a negative number and hence we use IG in Equation 18. We do not know the right value for τ^2 and hence we sample this parameter in a similar fashion as θ during the MCMC sampling process. We note that the prior for τ^2 must represent a distribution that can only sample positive real values, and we use the conjugate inverse Gamma prior with hyperparameters ν_1 and ν_2 representing the shape and scale parameters (see Section -B2).

We use the multivariate Gaussian distribution to represent the prior for parameters θ such as weights and bias of the linear model, which features negative and positive real numbers. Our model features more than one parameter, hence the multivariate Gaussian distribution is most appropriate for the prior. In this example, we adopt uninformative priors with hyperparameter values of $\sigma = 5$, $\nu_1 = 0$, and $\nu_2 = 0$ (Listing 6: Lines 15-17), but these values are user defined and could

be refined using trial runs. These values are based on expert opinion from analysis of related trained models.

First, we revisit multivariate normal distribution from Equation 4 to define the prior distribution for our linear model's parameters (weights and biases). Suppose that θ is our set of M parameters given by $(\theta = \theta_1, \dots, \theta_M)$. Since our prior is based on the normal distribution, we use mean $\mu = 0$ for each parameter to ensure we sample both positive and negative real numbers. Therefore, the mean μ is a vector of zeros and we get the prior using

$$f(\theta) = \frac{1}{\sqrt{(2\pi)^M |\Sigma|}} \exp \left(-\frac{1}{2} (\theta)^T \Sigma^{-1} (\theta) \right). \quad (19)$$

The covariance matrix Σ is a diagonal matrix with all values equal to σ^2 (scalar). Note that Σ^{-1} becomes I/σ^2 where I is an identity matrix (diagonal elements which are all ones). Hence, we take the numerator from Equation 19, i.e.

$$(\theta)^T \Sigma^{-1} (\theta) \quad (20)$$

becomes

$$\frac{(\theta)^T \mathbf{I} (\theta)}{\sigma^2}. \quad (21)$$

We note that multiplying identity matrix with any other matrix is the matrix itself, hence finally we get θ^2 in numerator. We can now move to the inverse-Gamma distribution used to define the prior for our model's variance (τ^2) and sampled (just as θ) and given by

$$f(\tau^2) = \frac{\nu_1^{\nu_2}}{\Gamma(\nu_1)} \left(\frac{1}{\tau^2} \right)^{\nu_1+1} \exp \left(\frac{-\nu_2}{\tau^2} \right). \quad (22)$$

We note that $\nu_1^{\nu_2}/\Gamma(\nu_1)$ is a constant which can be dropped considering proportionality. We take into account the product of all our sampled parameters (θ and τ^2) to define the combined prior, as given by Equation 23

$$p(\theta) \propto \frac{1}{(2\pi\sigma^2)^{M/2}} \times \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_{i=1}^M \theta^2 \right) \right\} \times \tau^{-2(1+\nu_1)} \exp \left(\frac{-\nu_2}{\tau^2} \right). \quad (23)$$

3) Python Implementation

The Python code presented in Listing 3³ begins the implementation of a Bayesian linear model using MCMC sampling. First, we define our simple linear model as given in Equation 12 using class LinearModel (Line 1) of Listing 3

³<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/03-Linear-Model.ipynb>

and define functions to evaluate the proposal (line 13), get the prediction (Line 25) and encode the parameters proposed from the MCMC sampler class into the linear model (Line 30).

Now that we have a class for our linear model, we can define the functions that will allow us to carry out MCMC sampling for the model parameters. We define our log-likelihood function using Equation 16, which becomes

$$\log p(\mathbf{y} | \mathbf{x}, \theta, \tau^2) = -\log((2\pi\tau^2)^{S/2}) - \frac{1}{2\tau^2} \sum_{t=1}^S (\mathbf{y}_t - f(\mathbf{x}_t, \theta))^2. \quad (24)$$

Furthermore, we define our log-prior using Equation 23 which becomes

$$\begin{aligned} \log p(\theta) \propto & -\frac{M}{2} \log 2\pi\sigma^2 \\ & -\frac{1}{2\sigma^2} \left(\sum_{i=1}^M \theta^2 \right) \\ & - (1 + \nu_1) \log \tau^2 - \frac{\nu_2}{\tau^2}. \end{aligned} \quad (25)$$

We present Python implementation of the log-likelihood (Lines 2 - 18) and the prior (Lines 21 - 37) in Listing 4.

Before running the MCMC sampler, we need to set up the sampler hyperparameters such as the maximum sampling time and burn-in period (Listing 7). We also need to assign hyperparameters that define the priors such as Gaussian prior variance (σ^2) and the IG prior parameters, ν_1 and ν_2 (Listing 6: Lines 15, 16 and 17). First, we need to generate an initial sample for our parameters, and initialise arrays to capture the samples that form the posterior distribution, the accuracy, and the model predictions as shown in code Listing 5 (Lines 5-23). Then we proceed with sampling as per the MCMC sampling algorithm detailed in Algorithm 1 and code Listing 5. This algorithm uses a Gaussian random walk for the parameter proposals (θ_p and τ_p^2), perturbing the previous proposed value with Gaussian noise as shown in Equations 26 and 27, respectively.

$$\theta_p \sim \theta_{p-1} + \mathcal{N}(0, \Delta_\theta) \quad (26)$$

$$\eta_p \sim \eta_{p-1} + \mathcal{N}(0, \Delta_\eta) \quad (27)$$

We implement the MCMC sampler with a Gaussian random-walk proposal for $\eta_p = \log \tau_p^2$, where we use η to represent τ^2 in log-space (Listing 5: Line 29). The step sizes for the proposals are determined by the hyperparameters Δ_θ and Δ_η which define the variance for the proposal of θ_p and η_p respectively. Once we sample η , we take the exponential to convert it back to the original form (see Line 30) and obtain τ^2 .

After getting the proposal for the parameters i.e. θ and τ^2 , we call the log-likelihood and prior functions to obtain

```

1 class LinearModel:
2     """
3     Simple linear model with a single output (y) given the covariates x_1...x_M of the form:
4     y = w_1 * x_1 + ... + w_M * x_M + b
5     where M = number of features, w are the weights, and b is the bias.
6     """
7     # Initialise values of model parameters
8     def __init__(self):
9         self.w = None
10        self.b = None
11
12    # Function to take in data and parameter sample and return the prediction
13    def evaluate_proposal(self, data, theta):
14        """
15        Encode the proposed parameters and then use the model to predict
16        Input:
17            data: (N x M) array of data
18            theta: (M + 1) vector of parameters. The last element of theta constitutes the bias term (
19                giving M + 1 elements)
20        """
21        self.encode(theta) # method to encode w and b
22        prediction = self.predict(data) # predict and return
23        return prediction
24
25    # Linear model prediction
26    def predict(self, x_in):
27        y_out = x_in.dot(self.w) + self.b
28        return y_out
29
30    # Helper function to split the parameter vector into w and band store in the model
31    def encode(self, theta):
32        self.w = theta[0:-1]
33        self.b = theta[-1]

```

Listing 3: Python implementation of a simple linear regression model

their respective values, as shown in Lines 32 - 35 of Listing 5. Note that the log-likelihood is used and hence the ratio of previous and current likelihood will need to consider log laws (rules), i.e., we note the log product rule in Equation 28 and the quotient rule in Equation 29. We use these rules in Lines 37 and 38 of Listing 5. Based on Equation 9 and taking the quotient rule into account since we are in the log space, we then accept/reject the proposed value according to the Metropolis-Hastings acceptance ratio (Line 42) as shown in Lines 41 - 54.

$$\log_b(x \times y) = \log_b(x) + \log_b(y) \quad (28)$$

$$\log_b(x/y) = \log_b(x) - \log_b(y) \quad (29)$$

Now that we have the sampler code (Listing 5), we can create an MCMC class that brings together the model, data, hyperparameters and sampling algorithm as shown in Listing 6.

E. BAYESIAN NEURAL NETWORKS VIA MCMC

1) Neural networks

We utilise a simple neural network, also known as a multilayer perceptron to demonstrate the process of training a Bayesian neural network via the MCMC sampler. A neural

network model $f(x)$ is made up of a series of lower level computations which can be used to transform inputs to their corresponding outputs $\{\bar{x}_t, y_t\}$. Neural networks feature layers of neurons whose value is determined based on a linear combination of inputs from the previous layer, with an activation function used to introduce nonlinearity.

We consider a simple neural network with one hidden layer with four input neurons, five hidden neurons and one output neuron, as shown in Figure 6. As an example, we can calculate the output value of the j th neuron in the first hidden layer of a network (h, j) using a weighted combination of the m inputs (\bar{x}_t), as shown in Equation 30.

$$g\left(\delta_{h,j} + \sum_{i=1}^m w_{i,j} \bar{x}_{t,i}\right) \quad (30)$$

where, the bias ($\delta_{h,j}$) and weights (w_i for each of the m inputs) are parameters to be sampled (trained or estimated), and $g(\cdot)$ is the activation function that is used to perform a nonlinear transformation. In our case, the function $g(\cdot)$ is the *sigmoid activation function*, used for the hidden and output layers as shown in Figure 6.

We train the model to approximate the function f such that $f(\mathbf{x}) = \bar{y}$ for all input-output pairs of instances from the training dataset. We extend our previous calculation, to a single neuron in the hidden layer to calculate the output $f(\bar{x}_t)$ as shown in Equation 31.

```

1 # Define the log-likelihood function
2 def likelihood_function(self, theta, tausq):
3     """
4     Calculate the likelihood of the data given the parameters
5     Input:
6         theta: (M + 1) vector of parameters. The last element of theta consitutes the bias term (giving M
7             + 1 elements)
8         tausq: variance of the error term
9     Output:
10         log_likelihood: log likelihood of the data given the parameters
11         model_prediction: prediction of the model given the parameters
12         accuracy: accuracy (RMSE) of the model given the parameters
13     """
14     # first make a prediction with parameters theta
15     model_prediction = self.model.evaluate_proposal(self.x_data, theta)
16     accuracy = self.rmse(model_prediction, self.y_data) #RMSE error metric
17     # now calculate the log likelihood
18     log_likelihood = np.sum(-0.5 * np.log(2 * np.pi * tausq) - 0.5 * np.square(self.y_data -
19         model_prediction) / tausq)
20     return [log_likelihood, model_prediction, accuracy]
21
22 # Define the prior
23 def prior(self, sigma_squared, nu_1, nu_2, theta, tausq):
24     """
25     Calculate the prior of the parameters
26     Input:
27         sigma_squared: variance of normal prior for theta
28         nu_1: parameter nu_1 of the inverse gamma prior for tau^2
29         nu_2: parameter nu_2 of the inverse gamma prior for tau^2
30         theta: (M + 1) vector of parameters. The last element of theta consitutes the bias term (giving M
31             + 1 elements)
32         tausq: variance of the error term
33     Output:
34         log_prior: log prior
35     """
36     n_params = self.theta_size # number of parameters in model
37     part1 = -1 * (n_params / 2) * np.log(sigma_squared)
38     part2 = 1 / (2 * sigma_squared) * (sum(np.square(theta)))
39     log_prior = part1 - part2 - (1 + nu_1) * np.log(tausq) - (nu_2 / tausq)
40     return log_prior

```

Listing 4: Python implementation of likelihood and prior functions for linear regression model to be incorporated into the MCMC sampling class

$$f(\mathbf{x}) = g\left(\delta_o + \sum_{j=1}^H v_h \times g\left(\delta_h + \sum_{i=1}^m w_{i,j} \mathbf{x}_i\right)\right) \quad (31)$$

where H is the number of neurons in the hidden layer, δ_o is the bias for the output, and v_h are the weights from the hidden layer to the output neuron. The complete set of parameters for the neural network model (Figure 6) is made up of $\theta = (\tilde{\mathbf{w}}, \tilde{\mathbf{v}}, \delta)$, where $\delta = (\delta_o, \delta_h)$. $\tilde{\mathbf{w}}$ are the weights transforming the input to hidden layer. $\tilde{\mathbf{v}}$ are the weights transforming the hidden to output layer. δ_h is the bias for the hidden layer, and δ_o is the bias for the output layer.

2) Bayesian neural networks

A Bayesian neural network is a probabilistic implementation of a standard neural network with the key difference being that the weights and biases are represented via the posterior probability distributions rather than single point values as shown in Figure 7. Similar to canonical neural networks [111], Bayesian neural networks also have universal continu-

ous function approximation capabilities. However, the posterior distribution of the network parameters allows uncertainty quantification on the predictions.

The task for MCMC sampling is to estimate (sample) the posterior distributions representing the weights and biases of the neural network that best fit the data. Perhaps, it can be argued that the method should be called an estimator, but we will stick to the sampler as given in the literature. As in the previous examples, we begin inference with prior distributions over the weights and biases of the network and use a sampling scheme to find the posterior distributions given training data. Since non-linear activation functions exist in the network, the conjugacy of prior and posterior is lost, and therefore we must employ an MCMC sampling scheme and make assumptions about the distribution of errors.

We specify the model similar to the Bayesian linear regression, assuming a Gaussian error as given in Equation 32.

$$y = f(\mathbf{x}, \theta) + e \quad e \sim \mathcal{N}(0, \tau^2) \quad (32)$$

This leads to the same likelihood function as presented in

```

1 # MCMC sampler
2 def sampler(self):
3     # Run the sampler for a defined linear model
4     # Define empty arrays to store the sampled posterior values
5     pos_theta = np.ones((self.n_samples, self.theta_size))
6     # posterior defining the variance of the noise in predictions
7     pos_tau = np.ones((self.n_samples, 1))
8     # record output f(x) over all samples
9     pred_y = np.ones((self.n_samples, self.x_data.shape[0]))
10    # record the RMSE of each sample
11    rmse_data = np.zeros(self.n_samples)
12
13    ## Initialisation: initialise theta - the model parameters
14    theta = np.random.randn(self.theta_size)
15    # make initial prediction
16    pred_y[0,] = self.model.evaluate_proposal(self.x_data, theta)
17    # initialise eta - we sample eta as a gaussian random walk in the log space of tau^2
18    eta = np.log(np.var(pred_y[0,] - self.y_data))
19    tausq_proposal = np.exp(eta)
20    # calculate the prior
21    prior_val = self.prior(self.sigma_squared, self.nu_1, self.nu_2, theta, tausq_proposal)
22    # calculate the likelihood considering observations
23    [likelihood, pred_y[0,], _, rmse_data[0]] = self.likelihood_function(theta, tausq_proposal)
24
25    ## Run the MCMC sample for n_samples
26    for ii in np.arange(1, self.n_samples):
27        # Sample new values for theta and tau using a Gaussian random walk
28        theta_proposal = theta + np.random.normal(0, self.step_theta, self.theta_size)
29        eta_proposal = eta + np.random.normal(0, self.step_eta, 1) # sample tau^2 in log space
30        tausq_proposal = np.exp(eta_proposal)
31        # calculate the prior
32        prior_proposal = self.prior(
33            self.sigma_squared, self.nu_1, self.nu_2, theta_proposal, tausq_proposal)
34        # calculate the log-likelihood considering observations
35        [likelihood_proposal, pred_y[ii,], _, rmse_data[ii]] = self.likelihood_function(theta_proposal,
36            tausq_proposal)
37        # Noting that likelihood_function and prior_val return log likelihoods, we can calculate the
38        # acceptance probability
39        diff_likelihood = likelihood_proposal - likelihood
40        diff_priorlikelihood = prior_proposal - prior_val
41        mh_prob = min(1, np.exp(diff_likelihood + diff_priorlikelihood))
42        # sample to accept or reject the proposal according to the acceptance probability
43        u = np.random.uniform(0, 1)
44        if u < mh_prob:
45            # accept and update the values
46            likelihood = likelihood_proposal
47            prior_val = prior_proposal
48            theta = theta_proposal
49            eta = eta_proposal
50            # store to make up the posterior
51            pos_theta[ii,] = theta_proposal
52            pos_tau[ii,] = tausq_proposal
53        else:
54            # reject move and store the old values
55            pos_theta[ii,] = pos_theta[ii-1,]
56            pos_tau[ii,] = pos_tau[ii-1,]
57    # store the posterior (samples after burn in) in a pandas dataframe and return
58    self.pos_theta = pos_theta[self.n_burnin:, ]
59    self.pos_tau = pos_tau[self.n_burnin:, ]
60    self.rmse_data = rmse_data[self.n_burnin:]
61    # split theta into w and b
62    results_dict = {'w{}'.format(_): self.pos_theta[:, _].squeeze() for _ in range(self.theta_size-1)}
63    results_dict['b'] = self.pos_theta[:, -1].squeeze()
64    results_dict['tau'] = self.pos_tau.squeeze()
65    results_dict['rmse'] = self.rmse_data.squeeze()
66    results_df = pd.DataFrame.from_dict(results_dict)
67    return results_df

```

Listing 5: Python implementation of an MCMC sampler for the linear model


```

1 class MCMC:
2     def __init__(self, n_samples, n_burnin, x_data, y_data):
3         self.n_samples = n_samples # number of MCMC samples
4         self.n_burnin = n_burnin # number of burn-in samples
5         self.x_data = x_data # (N x M)
6         self.y_data = y_data # (N x 1)
7         self.theta_size = x_data.shape[1] + 1 # weights for each feature and a bias term (M+1)
8
9         # MCMC sampler hyperparameters - defines the variance term in our Gaussian random walk
10        self.step_theta = 0.02;
11        self.step_eta = 0.01; # note eta is used as tau in the sampler to consider log scale.
12
13        # model hyperparameters
14        # considered by looking at distribution of similar trained models - i.e distribution of weights
15        # and bias
16        self.sigma_squared = 5
17        self.nu_1 = 0
18        self.nu_2 = 0
19
20        # initialise the linear model class
21        self.model = LinearModel()
22
23        # store output
24        self.pos_theta = None
25        self.pos_tau = None
26        self.rmse_data = None
27
28        # functions defined above - this is poor practice, but done for readability
29        # and clarity
30        self.likelihood_function = MethodType(likelihood_function, self)
31        self.prior_val = MethodType(prior_val, self)
32        self.sampler = MethodType(sampler, self)
33
34    def model_draws(self, num_samples = 10):
35        """
36        Simulate new model predictions (mu) under the assumption that our posteriors are
37        Gaussian.
38        """
39        # num_samples x num_data_points
40        pred_y = np.zeros((num_samples, self.x_data.shape[0]))
41        sim_y = np.zeros((num_samples, self.x_data.shape[0]))
42
43        for ii in range(num_samples):
44            theta_drawn = np.random.normal(self.pos_theta.mean(axis=0), self.pos_theta.std(axis=0), self.
45            theta_size)
46            tausq_drawn = np.random.normal(self.pos_tau.mean(), self.pos_tau.std())
47
48            [, pred_y[ii,:], sim_y[ii,:],_] = self.likelihood_function(
49                theta_drawn, tausq_drawn
50            )
51        return pred_y, sim_y

```

Listing 6: Python implementation of an MCMC class for Bayesian linear model.

logarithmic form in Equation 24. As in Section -D2, we adopt Gaussian priors for all parameters of the model (θ), with zero mean and a user-defined variance (σ^2), and an IG distribution for the variance of the error model (τ^2), with parameters ν_1 and ν_2 . The *likelihood function* and *prior* function, therefore remain unchanged from their definition in Listing 4.

3) Multinomial likelihood for classification problems

We note that neural network models are also prominent for classification problems apart from regression and prediction problems. Bayesian neural networks via the MCMC sampler require an appropriate likelihood function that is more suitable for discrete data to capture classification problems.

Hence, we use the multinomial likelihood, which is applicable to both binary and multi-class classification problems. Hence, we define the multinomial log-likelihood function for the classification problems using

$$\log(p(y|\theta)) = \sum_{t \in N} \sum_{k=1}^K z_{t,k} \log \pi_k \quad (33)$$

for classes $k = 1, \dots, K$, where π_k is the output of the neural network model after applying the transfer function, and N is the number of instances in the training data. In this case, we

```

1 ## MCMC Settings and Setup
2 n_samples = 20000 # number of samples to draw from the posterior
3 burn_in = int(n_samples* 0.25) # number of samples to discard before recording draws from the posterior
4
5 # Generate toy data
6 n_data = 100
7 n_features = 1
8 x_data = np.repeat(np.expand_dims(np.linspace(0, 1, n_data),axis=-1),n_features,axis=1)
9 y_data = 3 * x_data[:,0] + 4 + np.random.randn(n_data) * 0.5
10
11 # Initialise the MCMC class
12 mcmc = MCMC(n_samples, burn_in, x_data, y_data)
13 # Run the sampler
14 results = mcmc.sampler()
15 # Draw sample models from the posterior
16 pred_y, _ = mcmc.model_draws(num_samples=100)

```

Listing 7: Code to call the MCMC sampling class and fit a model to some toy data

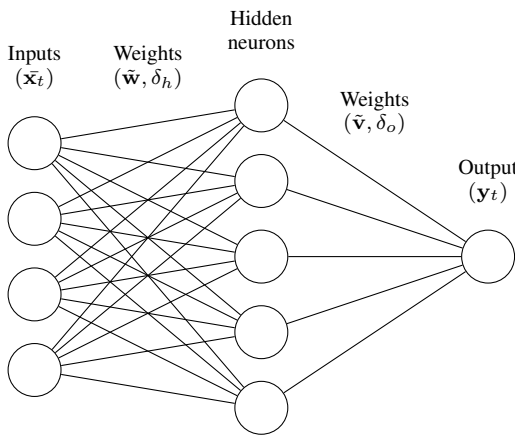


FIGURE 6: Simple neural network with a single hidden layer. The information is passed and processed from the input to hidden and then finally to the output layer.

utilize the *softmax function* [112] as the transfer function:

$$\pi_k = \frac{\exp(f(x_p))}{\sum_{k=1}^K \exp(f(x_k))} \quad (34)$$

for $k = 1, \dots, K$. $z_{t,k}$ is an indicator variable for the given instance of data t . We define class k in the data by

$$z_{t,k} = \begin{cases} 1, & \text{if } y_t = k \\ 0, & \text{otherwise.} \end{cases} \quad (35)$$

We note that we do not use the noise parameter (i.e., τ^2) as in the case of the inverse gamma distribution for the Gaussian likelihood earlier, hence we do not need to have a prior distribution for this case. We will use a Gaussian prior for weights and biases of the neural network model. Therefore, in the case of classification, our prior distribution from Equation 25 simplifies to

$$p(\theta) \propto \frac{1}{(2\pi\sigma^2)^{M/2}} \times \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_{i=1}^M \theta^2 \right) \right\}. \quad (36)$$

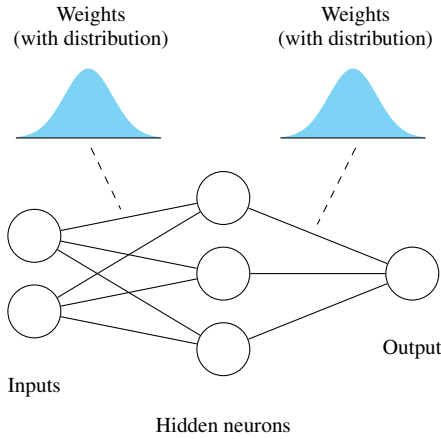
and finally log-prior becomes

$$\log p(\theta) \propto -\frac{M}{2} \log 2\pi\sigma^2 \times -\frac{1}{2\sigma^2} \left(\sum_{i=1}^M \theta^2 \right) \quad (37)$$

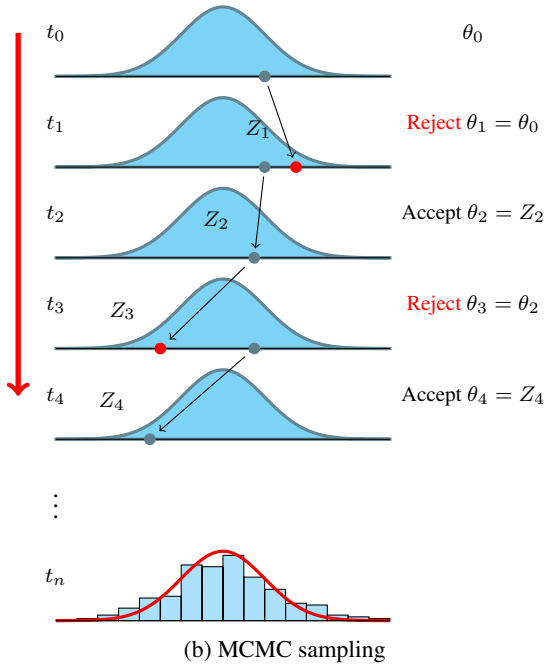
4) Training neural networks via backpropagation

We note that typically random-walk proposal distributions are used for small scale-models such as linear models; however, neural network models features a large number of parameters. The choice of a proposal distribution is essential for models with large number of parameters. We need to incorporate gradients into our proposal distribution for better sampling, and we will begin by examining how gradients are incorporated in conventional neural networks and deep learning models.

Gradient-based methods have been widely used in machine learning and served as the backbone in the backpropagation algorithm [113]. A prominent implementation is *stochastic gradient descent* (SGD) which involves stepping through the parameter space iteratively in a stochastic manner using gradients, to optimize a differentiable objective function. The method has been prominently featured in the backpropagation algorithm for training of various neural networks architectures including deep learning models [114]–[116]. Backpropagation involves a forward pass which propagates information forward to get prediction (decision) at the output layer and a backward pass to compute the local gradients for each of the parameters. These gradients are then used to inform the update of the model parameters in an iterative process, where the weights and biases are updated at each



(a) Bayesian neural network



(b) MCMC sampling

FIGURE 7: Bayesian neural network and MCMC sampling adapted from [104].

step. The training of neural networks is also considered as solving a non-convex optimization problem $\text{argmin}_{\theta} L(\theta)$; where, $\theta \in R^n$ is the set of parameters and L is the loss function. We give the parameter (weight) update for an iteration (epoch) of SGD in Equation 38

$$\theta_k = \theta_{k-1} - a_k \nabla L(\theta_{k-1}) \quad (38)$$

where, θ_k denotes the k^{th} iteration, a_k is the learning rate, and $\nabla L(\theta_k)$ denotes the gradient.

We note that the *learning rate* is user defined hyperparameter which depends on the problem and data at hand. It is typically determined through tuning using cross-validation or trial and error. A number of extensions of the backpropagation algorithm employing SGD have been proposed to

address limitations. These include use of *weight decay* regularization during training to improve generalisation ability [117], a momentum mechanism for faster training [118], [119], adaptive learning rate [119], and second-order gradient methods [120] which, although efficient, have problems in scaling up computationally with larger models. In the last decade, with the deep learning revolution, further attempts have been made to ensure that enhanced backpropagation algorithms, not only improve training accuracy, but can also scale well computationally for large deep learning models. These led to the development of methods such as the adaptive gradient algorithm (AdaGrad) [121], Ada-Delta [122], and Adam (adaptive moment estimation) [123]. These advanced algorithms are generally based on the idea of adapting the learning rate automatically during the training, taking into account the recent history of the optimisation process.

a: Langevin proposal distribution

We mentioned earlier that random-walk proposal distributions are suited for small scale models and better proposal distributions would be required for neural network models. Although simple neural networks have much lower number of parameters, when compared to deep learning models, training simple neural networks with MCMC sampling is a challenge with random-walk proposal distribution. We need to utilise the properties of backpropagation algorithm and the mechanism of weight update using gradients. Hence, we utilise *stochastic gradient Langevin dynamics* [61] for the proposal distribution which features the addition of noise to the stochastic gradients. The method has shown to be effective for linear models [61] which motivated its use in Bayesian neural networks. In the literature, Langevin MCMC has been very promising for simple and deep neural networks [62], [69], [70]. Hence, we draw the proposed values for the parameters (θ^p) according to a one-step (epoch) gradient as shown in Equation 39.

$$\theta^p \sim \mathcal{N}(\bar{\theta}^{[s]}, \Sigma_{\theta}) \quad (39)$$

A Gaussian distribution with a standard deviation of Σ_{θ} , and mean ($\bar{\theta}^{[s]}$) calculated using a gradient based update (Equation 40 of the parameter values from the previous step ($\theta^{[s]}$).

$$\bar{\theta}^{[s]} = \theta^{[s]} + r \times \nabla E(\theta^{[s]}) \quad (40)$$

with learning rate r and gradient update ($\nabla E(\theta^{[s]})$) according to the model residuals.

$$E(\theta^{[s]}) = \sum_{t \in \mathcal{T}} (y_t - F(x_i, \theta^{[s]}))^2 \quad (41)$$

$$\nabla E(\theta^{[s]}) = \left(\frac{\partial E}{\partial \theta_1}, \dots, \frac{\partial E}{\partial \theta_L} \right).$$

Hence, the Langevin proposal distribution (also known as Langevin-gradient) consists of 2 parts:

- 1) Gradient descent-based weight update

2) Addition of Gaussian noise from $\mathcal{N}(0, \Sigma_\theta)$

We need to ensure that the detailed balance is maintained since the Langevin-gradient proposals are not symmetric. We note that MCMC implementations with relaxed detailed balance conditions for some applications also exist [124]. Therefore, we use a combined update in the Metropolis-Hastings step, which accepts the proposal θ^p for a position s with the probability α , as shown in Equation 42

$$\alpha = \min \left\{ 1, \frac{p(\theta^p|\mathbf{y})q(\theta^{[s]}|\theta^p)}{p(\theta^{[s]}|\mathbf{y})q(\theta^p|\theta^{[s]})} \right\} \quad (42)$$

where $p(\theta^p|\mathbf{y})$ and $p(\theta^{[s]}|\mathbf{y})$ can be computed using the likelihood and prior (Equations (16) and (23)). We give the ratio of the proposed and the current $q(\theta^p|\theta^{[s]})$ in Equation 43

$$q(\theta^{[s]}|\theta^p) \sim N(\bar{\theta}^{[s]}, \Sigma_\theta) \quad (43)$$

which is based on a one-step (epoch) gradient $\nabla E_y[\theta^{[s]}]$ and learning rate r , as given in Equation 44

$$\bar{\theta}^{[s]} = \theta^{[s]} + r \times \nabla E_y[\theta^{[s]}]. \quad (44)$$

Thus, this ensures that the detailed balance condition holds and the sequence $\theta^{[s]}$ converges to draws from the posterior $p(\theta|y)$. Since our implementation is in the log-scale, we give the log-posterior in Equation 45

$$\log(p(\theta|y)) = \log(p(\theta)) + \log(p(y|\theta)) + \log(q(\theta|\theta^*)) \quad (45)$$

Algorithm 2 gives a full description of the Langevin MCMC sampling scheme with user-defined parameters that include the maximum number of samples (S_{max}), rate of Langevin-gradient proposals (L_{prob}), and learning rate r used for the Langevin-gradient proposals. We note that in a standard Langevin MCMC approach, $L_{prob} = 1$ and Gaussian noise is already part of Langevin-gradient distribution. However, in our implementation, we use a combination of random-walk proposal distribution with Langevin-gradients, as this is computationally more efficient than Langevin-gradients alone, which consumes extra computational time when computing gradients, especially in larger models.

We begin by drawing initial values for the θ from the prior distribution given in Equation (23) (Stage 1.1). We draw a new proposal for θ^p (which incorporates the model weights and biases and τ^2 from either a Langevin-gradient or random-walk proposal distribution (Stage 1.2). We evaluate the proposal using the Bayesian neural network (BNN) model with the log-likelihood function in Equation 16 (Stage 1.4) and the prior given in Equation (23) (Stage 1.3). We can then check if the proposal should be accepted using Metropolis-Hastings condition (Stage 1.5 and 1.6). If accepted, the proposal becomes part of the chain, else we retain previous the last accepted state of the chain. We repeat the procedure until the maximum samples are reached (S_{max}).

Finally, we execute post-sampling stage, where we obtain the posterior distribution by concatenating the history of the samples in the chain.

Data: Dataset

Result: Posterior distribution of model parameters (weights and biases)

- Stage 1.0: Metropolis Transition

- 1.1 Draw initial values θ_0 from the prior

for each s until S_{max} do

1.2 Draw κ from a Uniform-distribution [0,1] **end**

if $\kappa \leq L_{prob}$ then

Use Langevin-gradient proposal distribution:

$$\theta^p \sim \mathcal{N}(\bar{\theta}^{[s]}, \Sigma_\theta)$$

end

else

Use random-walk proposal distribution:

$$\theta^p \sim \mathcal{N}(\theta^{[s]}, \Sigma_\theta)$$

end

1.3 Evaluate prior given in Equation 23

1.4 Evaluate log-likelihood given in Equation 16

1.5 Compute the posterior probability for

Metropolis-Hastings condition - Equation 45

$$\log(p(\theta|y)) = \log(p(\theta)) + \log(p(y|\theta)) + \log(q(\theta|\theta^*))$$

1.6 Draw u from a Uniform-distribution [0,1]

if $\log(u) \leq \log(p(y|\theta))$ then

Accept replica state: $\theta^{[s+1]} \leftarrow \theta^p$

end

else

Reject and retain previous state: $\theta^{[s+1]} \leftarrow \theta^{[s]}$

end

Algorithm 2: Bayesian neural network via Langevin MCMC sampling.

5) Python Implementation

We first define and implement the simple neural network module (class), and implement methods (functions) for the forward and backward pass to calculate the output of the network given a set of inputs. We need to compute the gradients and update the model parameters given a model prediction and observations, respectively. Listings 8, 9, and 11 presents the implementation⁴ of the Bayesian Neural Network and associated Langevin MCMC sampling scheme. Note that we implement the Bayesian Neural Network via the MCMC sampler class to sample (train) the weights and biases of the Neural Network class.

Next, we implement the model for a single hidden layer neural network with multiple input neurons and multiple output neurons (for binary and multi-class classification and multi-output regression). Listing 8 defines the Neural Network class with the *constructor* function (*init*) which defines the network topology, in terms of the number of

⁴<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/04-Bayesian-Neural-Network.ipynb>


```

1 class NeuralNetwork:
2     """
3     Neural Network model with a single hidden layer and a single output (y)
4     """
5     def __init__(self, layer_sizes, learning_rate=0.01):
6         """ Initialize the model
7         Input:
8             - layer_sizes (input, hidden, output): array specifying the number of
9               nodes in each layer
10            - learning_rate: learning rate for the gradient update
11         """
12        # Initial values of model parameters
13        self.input_num = layer_sizes[0]
14        self.hidden_num = layer_sizes[1]
15        self.output_num = layer_sizes[2]
16        # total number of parameters from weights and biases
17        self.n_params = (self.input_num * self.hidden_num) + (self.hidden_num * self.output_num) + \
18            self.hidden_num + self.output_num
19        # learning params
20        self.lrate = learning_rate
21        # Initialize network structure
22        self.initialise_network()
23        # functions defined above - this is poor practice, but done for readability
24        # and clarity
25        self.forward_pass = MethodType(forward_pass, self)
26        self.backward_pass = MethodType(backward_pass, self)
27
28    def initialise_network(self):
29        """
30        Initialize network structure - weights and biases for the hidden layer and output layer
31        """
32        # hidden layer
33        self.l1_weights = np.random.normal(
34            loc=0, scale=1/np.sqrt(self.input_num),
35            size=(self.input_num, self.hidden_num))
36        self.l1_biases = np.random.normal(
37            loc=0, scale=1/np.sqrt(self.hidden_num),
38            size=(self.hidden_num,))
39        # placeholder for storing the hidden layer values
40        self.l1_output = np.zeros((1, self.hidden_num))
41        # output layer
42        self.l2_weights = np.random.normal(
43            loc=0, scale=1/np.sqrt(self.hidden_num),
44            size=(self.hidden_num, self.output_num))
45        self.l2_biases = np.random.normal(
46            loc=0, scale=1/np.sqrt(self.hidden_num),
47            size=(self.output_num,))
48        # placeholder for storing the model outputs
49        self.l2_output = np.zeros((1, self.output_num))
50
51    def evaluate_proposal(self, x_data, theta):
52        """
53        A helper function to take the input data and proposed parameter sample and return the prediction
54        Input:    data: (N x num_features) array of data
55                theta: (w,v,b_h,b_o) vector of parameters with weights and biases
56        """
57        self.decode(theta) # method to decode w into W1, W2, B1, B2.
58        size = x_data.shape[0]
59        fx = np.zeros(size)
60        for i in range(0, size): # to see what fx is produced by your current weight update
61            fx[i] = self.forward_pass(x_data[i,])
62        return fx
63
64    def sigmoid(self, x):
65        # sigmoid activation function
66        return 1 / (1 + np.exp(-x))
67
68    def softmax(self, x):
69        #softmax function
70        prob = np.exp(x) / np.sum(np.exp(x))
71        return prob

```

Listing 8: Python implementation of the Neural Network class

```

1 # NN prediction
2 def forward_pass(self, X):
3     """
4     Take an input X and return the output of the network
5     Input:
6         - X: (N x num_features) array of input data
7     Output:
8         - self.l2_output: (N) array of output data f(x) which can be
9           compared to observations (Y)
10    """
11    # Hidden layer
12    l1_z = np.dot(X, self.l1_weights) + self.l1_biases
13    self.l1_output = self.sigmoid(l1_z) # activation function g(.)
14    # Output layer
15    l2_z = np.dot(self.l1_output, self.l2_weights) + self.l2_biases
16    self.l2_output = self.sigmoid(l2_z)
17    return self.l2_output
18
19 def backward_pass(self, X, Y):
20     """
21     Compute the gradients using a backward pass and undertake Langevin-gradient updating of parameters
22     Input:
23         - X: (N x num_features) array of input data
24         - Y: (N) array of target data
25    """
26    # dE/dtheta
27    l2_delta = (Y - self.l2_output) * (self.l2_output * (1 - self.l2_output))
28    l2_weights_delta = np.outer(
29        self.l1_output,
30        l2_delta
31    )
32    # backprop of l2_delta and same as above
33    l1_delta = np.dot(l2_delta, self.l2_weights.T) * (self.l1_output * (1 - self.l1_output))
34    l1_weights_delta = np.outer(
35        X,
36        l1_delta
37    )
38
39    # update for output layer
40    self.l2_weights += self.lrate * l2_weights_delta
41    self.l2_biases += self.lrate * l2_delta
42    # update for hidden layer
43    self.l1_weights += self.lrate * l1_weights_delta
44    self.l1_biases += self.lrate * l1_delta

```

Listing 9: Python implementation of Neural Network forward and backward passes

input, hidden and output neurons along with the learning rate. These values are passed by the calling function. Next, we compute the total number of parameters by Line 17. In Line 22, we initialize the network by calling the function (*initialise_network*) where we initialize (create) matrices for the weights from the input-hidden, and hidden-output layer, along with the vectors for their biases (Lines 33-49). Line 51 gives the *evaluate_proposal* function that takes the input data and proposed parameters (*x_data* and *theta*) and returns the prediction (*fx*) in Line 62. We feature the sigmoid transfer (activation) function in Line 64.

Listing 9 lists the rest of the functions from the Neural Network class, where *forward_pass* propagates the information forward from input - hidden layer and then hidden to output layer using a *dot product* (Lines 12 and 16) and the *returns* the output layer (Line 17). The *backward_pass* function begins by computing the gradients (delta) at the output layer (Line 27) and hidden layer (Line 33). Lines

40–44 update the weights in the hidden and output later using their respective gradients.

In a conventional backpropagation algorithm implementation, basically the forward and backward pass functions will be called in an iterative loop that will call these functions until the maximum number of epochs, or a given training (validation) error, has been reached. However, in our case, we are using the Langevin MCMC sampler to train the neural network model; hence, we have additional helper functions (Listing 10) to ensure that the MCMC sampler class gets the information as needed. Essentially, from the MCMC class (Listing 11), the sampler function (Listing 13) calls the respective functions to evaluate the likelihood and the prior (Listing 12). In the case of computing the likelihood, the *evaluate_proposal* function in Listing 8 calls the decode function in Listing 10 to insert the values of the proposal (*theta*) into the weight matrices and bias vectors of the model defined in the *NeuralNetwork* class of Listing 8.

```

1  def langevin_gradient(self, x_data, y_data, theta, depth):
2      """
3      Compute the Langevin-gradient proposal distribution
4      Input:
5          - x_data: (N x num_features) array of input data
6          - y_data: (N) array of target data
7          - theta: (w,v,b_h,b_o) vector of proposed parameters.
8          - depth: SGD depth
9      Output:
10         - theta_updated: Updated parameter proposal
11      """
12      self.decode(theta) # method to decode w into W1, W2, B1, B2.
13      size = x_data.shape[0]
14      # Update the parameters based on LG
15      for _ in range(0, depth):
16          for ii in range(0, size):
17              self.forward_pass(x_data[ii,:])
18              self.backward_pass(x_data[ii:], y_data[ii])
19      theta_updated = self.encode()
20      return theta_updated
21
22  def encode(self):
23      """
24      Encode the model parameters into a vector
25      theta: vector of parameters.
26      """
27      w1 = self.l1_weights.ravel()
28      w2 = self.l2_weights.ravel()
29      theta = np.concatenate([w1, w2, self.l1_biases, self.l2_biases])
30      return theta
31
32  def decode(self, theta):
33      """
34      Decode the model parameters from a vector
35      theta: vector of parameters.
36      """
37      w_layer1size = self.input_num * self.hidden_num
38      w_layer2size = self.hidden_num * self.output_num
39      w_layer1 = theta[0:w_layer1size]
40      self.l1_weights = np.reshape(w_layer1, (self.input_num, self.hidden_num))
41
42      w_layer2 = theta[w_layer1size:w_layer1size + w_layer2size]
43      self.l2_weights = np.reshape(w_layer2, (self.hidden_num, self.output_num))
44      self.l1_biases = theta[w_layer1size + w_layer2size:w_layer1size + w_layer2size + self.hidden_num]
45      self.l2_biases = theta[w_layer1size + w_layer2size + self.hidden_num:w_layer1size + w_layer2size
+ self.hidden_num + self.output_num]

```

Listing 10: Langevin-gradient functions in the Neural Network class

Next, we move to Listing 11 that implements the Langevin MCMC sampler for classification problems, as given in the notebook of Github repository ⁵. We note that we also provide the implementation for regression/prediction problems in the notebook ⁶. Furthermore, we also provide Python code implementation that features both classification and regression problems in the repository ⁷.

In Listing 11, we define the MCMC class with number of samples ($n_{samples}$), the burnin period (n_{burnin}), along with the training (x_{data} and y_{data}) and test datasets (x_{test}

⁵<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/04a-Bayesian-Neural-Network-Classification.ipynb>

⁶<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/04-Bayesian-Neural-Network.ipynb>

⁷<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/tree/main/code>

and y_{test}). Lines 12-14 initializes the hyperparameters, such as the $step_size$ of the random-walk on θ and the $\sigma_squared$ that defined the spread of the Gaussian prior for the weights and biases. Lines 17-21 define the neural network model, the use of Langevin-gradients, the probability (l_prob) for using it, and the total number of weights and biases (θ_size). Next comes the storage of the parameters that are samples (Lines 24-25). Line 27 defines the function for $model_draws$ - this is used post-sampling as a means to test the trained model. Line 50 defines the classification prediction accuracy, note that other error metrics can also be added.

Listing 12 defines the multinomial Log-likelihood and prior for classification problems. Line 1 implements the multinomial log-likelihood function that uses the parameters and data (Equation 33). Line 22 implements the log-prior function that uses the proposals (θ) and the user-defined

variance (*sigma_squared*) for the Gaussian prior. Note that in the case of regression and prediction problems, the log-likelihood and log-prior are similar to the Bayesian linear regression (Listing 4) with the omission of terms related to τ^2 .

In Listing 14, we begin sampling by first initializing variables that track the number of accepted proposals and how many times Langevin-gradients are utilized, this is just for analysis. In Line 4, we begin the sampling using a for loop that begins with 1 and ends with the number of samples. In Line 6, we propose the new values for the parameters (*theta*) using random-walk proposal distribution centred at the mean of 0 and given spread (*step_size*) which needs to be experimentally determined in trial experiments. Then, we decide if we wish to use the Langevin-gradients or random-walk proposal distribution (Lines 7-8). Lines 9-20 implement the Langevin-gradients where we get one-step gradients from the neural network model. Therefore, we need to run a forward-pass and a backward-pass using the new sets of the current weights and biases, obtain the gradients of the output and hidden layers of the network, and concatenate to return these as a vector (Lines 15-20) of Listing 10.

In Listing 14, we then use the gradient (*theta_grad*) as the centre for the normal distribution to draw and add Gaussian noise to the gradients (Line 10). Then we again obtain the gradients, this is merely for the detailed balance condition. In Line 11, we get the gradients again, but this time we use the new values of *theta* (*theta_proposal*) that we earlier obtained in Line 6. As given in Equation 42, in the case when the proposals are not symmetric (i.e., Langevin-gradients), we need to get the q-ratio (Line 19 *diff_prop*). In the log-scale, this is obtained by the difference in the current *theta* (first) and the new *theta_proposal* (second) to account for the detailed balance condition as shown in Line 19. In order to obtain the q-ratio, we need to further evaluate the old and the new proposals using the multivariate normal distribution and for numerical stability. However, we need to have a simplified implementation for the multivariate distribution given a large set of weights and biases (*theta*). Since we are operating in the log-scale, we can further simplify the multivariate normal distribution as shown in Lines 14-19. Finally, Line 23 implements the case when random-walk proposal distribution would be used, note that in Line 22, the *diff_prop* is 0. This accounts for the detailed balance condition, since the proposals are naturally symmetric, in the case of random-walk proposal distribution (Line 23).

Next, we compute the log-likelihood in Listing 14 (Line 27) and the error metrics for the test dataset (Line 29). We note that this is a classification problem, and hence the classification accuracy (Listing 12 - Line 13) is reported using the log-likelihood function. We determine the Metropolis-Hastings (MH) acceptance rate using Lines 34-35. In Lines 32 and 33, we get the difference (ratio) for the proposed likelihood and the current likelihood, and the ratio for the prior with the current and proposed value of the prior. We utilize these to get the MH probability in Line 34, which

also utilizes the difference (ratio) of proposed and current proposals, obtained either from Line 19 or Line 22. Finally, we either accept (Lines 39-45) or reject (Line 48) the proposal by comparing the MH probability with a random value obtained in Line 35. In the case if the proposal is accepted, the proposed values of *theta* along with prior and likelihood, becomes the current value in the chain. In the case if it is rejected, then the chain maintains its last accepted value as the current value (Line 48). We remove the burn-in portion and store the posterior (Lines 53-57). Finally, we return the dictionary of the data that features the posterior and predictions using the Pandas library.

F. RESULTS

We use the Sunspot time series⁸ data and Abalone⁹ datasets for regression problems. The Abalone dataset provides the ring age for Abalone based on eight features that represent physical properties such as length, width, and weight and associated target feature, i.e., the ring age. We note that determining the age of Abalone is difficult, as requires cutting the shell and counting the number of rings using a microscope. However, other physical measurements can be used to predict the age and a model can be developed to use the physical features to determine the ring age. Sunspots are regions of reduced surface temperature in the Sun's photosphere caused by concentrations of magnetic field flux, and appears as spots darker than the surrounding areas. The Sunspot cycles are about every eleven years and over the solar cycle, the number of Sunspot changes more rapidly. Sunspot activities are monitored since they have an include on Earth's climate and weather systems. We obtain the Abalone dataset from the University of California (UCI) Machine Learning Repository¹⁰ and keep a processed version of all the datasets in our repository¹¹.

In the Sunspot time series problem, we employ a one-step ahead prediction and hence use one output neuron. We process the Sunspot dataset (univariate time series) using Taken's embedding theorem [125] to construct a state-space vector, also known as *data windowing*. This is essentially using a sliding window approach of size D overlapping T time lags. The window size D determines the number of input neurons in the Bayesian neural network and Bayesian linear model. We used $D = 4$ and $T = 2$ for our data reconstruction for the Sunspot time series, as these values have given good performance in our previous works [62].

We also obtain datasets for classification problems from the same repository that features a large number of datasets for classification problems. We selected the Iris classification dataset that contains 4 features (sepal length, sepal width, petal length, petal with) of three types of Iris flower species,

⁸<https://www.sidc.be/silso/datafiles>

⁹<https://archive.ics.uci.edu/ml/datasets/abalone>

¹⁰<https://archive-beta.ics.uci.edu/about>

¹¹<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/tree/main/data>


```

1
2 class MCMC:
3     def __init__(self, model, n_samples, n_burnin, x_data, y_data, x_test, y_test):
4         self.n_samples = n_samples # number of MCMC samples
5         self.n_burnin = n_burnin # number of burn-in samples
6         self.x_data = x_data # (N x num_features)
7         self.y_data = y_data # (N x 1)
8         self.x_test = x_test # (Nt x num_features)
9         self.y_test = y_test # (Nt x 1)
10
11         # MCMC parameters - defines how much variation you need in changes to theta, tau
12         self.step_theta = 0.025;
13         # Hyperpriors
14         self.sigma_squared = 25
15
16         # initialise the neural network model class
17         self.model = model
18         self.use_langevin_gradients = True
19         self.sgd_depth = 1
20         self.l_prob = 0.5 # likelihood prob
21         self.theta_size = self.model.n_params # weights for each feature and a bias term
22
23         # store output
24         self.pos_theta = None
25         self.rmse_data = None
26
27     def model_draws(self, num_draws = 10, verbose=False):
28         ''' Calculate the output of the network from draws of the posterior distribution
29         Input: num_draws: number of draws, verbose: if True, print the details of each draw
30         Output: pred_y: (num_draws x N) ouptut of the NN for each draw'''
31         accuracy = np.zeros(num_draws)
32         pred_y = np.zeros((num_draws, self.x_data.shape[0]))
33         sim_y = np.zeros((num_draws, self.x_data.shape[0]))
34
35         for ii in range(num_draws):
36             theta_drawn = np.random.normal(self.pos_theta.mean(axis=0), self.pos_theta.std(axis=0), self.
theta_size)
37             [likelihood_proposal, pred_y[ii,], sim_y[ii,], accuracy[ii]] = self.likelihood_function(
38                 theta_drawn
39             )
40             if verbose:
41                 print(
42                     'Draw {} - accuracy: {:.3f}. Theta: {}'.format(
43                         ii, accuracy[ii], theta_drawn
44                     )
45                 )
46         return pred_y, sim_y
47
48     # Additional error metric
49     @staticmethod
50     def accuracy(predictions, targets):
51         '''
52         Additional error metric - accuracy
53         '''
54         count = (predictions == targets).sum()
55         return 100 * (count / predictions.shape[0])

```

Listing 11: Bayesian neural network using MCMC sampler for classification problems

```

1 def likelihood_function(self, theta, test=False):
2     ''' Calculate the multinomial log-likelihood of the data given the parameters and model
3     Input:  theta: vector of parameters
4     Output: log_likelihood: log likelihood of the data given the parameters '''
5     if test:
6         x_data = self.x_test
7         y_data = self.y_test
8     else:
9         x_data = self.x_data
10        y_data = self.y_data
11        model_prediction, probs = self.model.evaluate_proposal(x_data, theta)
12        model_simulation = model_prediction
13        accuracy = self.accuracy(model_prediction, y_data) #Accuracy error metric
14        # now calculate the log-likelihood
15        log_likelihood = 0
16        for ii in np.arange(x_data.shape[0]):
17            for jj in np.arange(self.model.output_num):
18                if y_data[ii] == jj:
19                    log_likelihood += np.log(probs[ii, jj])
20        return [log_likelihood, model_prediction, model_simulation, accuracy]
21
22 def prior(self, sigma_squared, theta):
23     '''
24     Calculate the prior of the parameters
25     Input: sigma_squared - variance of normal prior for theta
26     Output: log_prior
27     '''
28     n_params = self.theta_size # number of parameters in model
29     part1 = -1 * (n_params / 2) * np.log(sigma_squared)
30     part2 = 1 / (2 * sigma_squared) * (sum(np.square(theta)))
31     log_prior = part1 - part2
32     return log_prior

```

Listing 12: Multinomial Log-likelihood and prior for classification problems (Continued from Listing 11)

featuring 50 instances for each case¹². This dataset is one of the most prominent datasets used for machine learning. We also selected the Ionosphere dataset that features a binary classification task with 351 instances¹³. It has 34 continuous features and the task is to filter the radio signals as "good" or "bad".

In the Bayesian linear model and neural network, we choose the number of samples to be 25,000 for all problems, distributed across 5 chains and excluding 50% burn-in. In the Bayesian linear model, we choose the learning rate $r = 0.1$, and the step sizes for $\theta = 0.02$ and $\tau = 0.01$, respectively. Additionally, for the Gaussian prior distribution, we choose the parameters $\sigma^2 = 5$, $\nu_1 = 0$ and $\nu_2 = 0$, respectively. In the Bayesian neural network models, we choose the learning rate $r = 0.01$, and the step sizes for $\theta = 0.025$ and $\tau = 0.2$, respectively. In the Gaussian prior distribution, we choose the hyperparameter $\sigma^2 = 25$ determined from examining trained neural network models for similar problems. In the case of regression, we use inverse-Gamma prior for τ^2 , and hence use hyperparameters for this prior, $\nu_1 = 0$ and $\nu_2 = 0$, respectively. We also use a burn-in rate of 0.5 for both the Bayesian linear and the Bayesian neural network models. In the case of Bayesian neural networks, we apply Langevin-gradients at a rate of 0.5.

We first present the results of Bayesian regression with

the Sunspot (time series prediction) and Abalone (regression) datasets. We evaluate the model performance using the *root mean squared error* (RMSE) which is a standard metric for time series prediction and regression problems. We present the results obtained by the Bayesian linear model and Bayesian neural network model for regression problems in Table 1. Figures 8 and 9 present the prediction plots (observed, modelled and 95 % credible interval) that show a comparison between Bayesian linear model and neural network models for the fixed training (Panels a and b) and test (Panels c and d) datasets. We note that Panels a and c show the timeseries prediction (x-axis represents timestep number). Panels b and d present a scatter-plot of the change from timestep $t - 1$ to t in the observed (ΔY observed) and predicted values (ΔY modelled). This gives an indication of model's ability to predict change at each timestep with a skill better than persistence (as observed Y_{t-1} is given as an input to the model, a model predicting $y_t = Y_{t-1}$ could have a low RMSE). Thus, we note that the RMSE does not assess the skill (capability) of the model above the observed Y_{t-1} (given these models are conducting one step-ahead prediction); however, given the analysis using Panels b and d, the performance of the two models can be further compared using RMSE.

In Table 1, we observe that the Bayesian neural network performs better for the Sunspot time series prediction problem, as it achieves a better accuracy (lower RMSE) on both

¹²<https://archive.ics.uci.edu/ml/datasets/iris>

¹³<https://archive.ics.uci.edu/ml/datasets/ionosphere>

```

1  # MCMC sampler
2  def sampler(self):
3      """
4      Run the sampler for a defined Neural Network model
5      """
6      # define empty arrays to store the sampled posterior values
7      # posterior of all weights and bias over all samples
8      pos_theta = np.ones((self.n_samples, self.theta_size))
9
10     # record output f(x) over all samples
11     pred_y = np.zeros((self.n_samples, self.x_data.shape[0]))
12     # record simulated values f(x) + error over all samples
13     sim_y = np.zeros((self.n_samples, self.x_data.shape[0]))
14     # record the RMSE of each sample
15     accuracy_data = np.zeros(self.n_samples)
16     # now for test
17     test_pred_y = np.ones((self.n_samples, self.x_test.shape[0]))
18     test_sim_y = np.ones((self.n_samples, self.x_test.shape[0]))
19     test_accuracy_data = np.zeros(self.n_samples)
20
21     ## Initialisation
22     # initialise theta
23     theta = np.random.randn(self.theta_size)
24     # make initial prediction
25     pred_y[0,], _ = self.model.evaluate_proposal(self.x_data, theta)
26
27     # Hyperparameters of priors - considered by looking at distribution of similar trained models - i.e
28     # distribution of weights and bias
29     sigma_squared = self.sigma_squared
30
31     # calculate the prior
32     prior_val = self.prior(sigma_squared, theta)
33     # calculate the likelihood considering observations
34     [likelihood, pred_y[0,], sim_y[0,], accuracy_data[0]] = self.likelihood_function(theta)

```

Listing 13: Implementation of MCMC sampler function (Continued from Listing 12)

the training and testing set. This can also be seen in Figures 8 and 9. In the case of the Abalone problem, Table 1 shows that both models obtain similar classification performance, but the Bayesian neural network has better test performance. However, we note that in both problems, Bayesian neural networks have a much lower acceptance rate; we prefer roughly a 23 % acceptance rate [126] that implies that the posterior distribution has been effectively sampled. However, we also note that such MCMC sampling acceptance rates are typically based on statistical and linear models, which may not apply to Bayesian neural networks. Therefore, more research needs to be done to determine a good acceptance rate that aligns with convergence and *ergodicity* [127].

Table 2 presents results for the classification problems in the Iris and Ionosphere datasets. We notice that both models have similar test and training classification performance for the Iris classification problem, and Bayesian neural networks give better results for the test dataset for the Ionosphere problem. The acceptance rate is much higher for Bayesian neural networks in the case of classification, and in the case of regression, only the Bayesian linear model has a more suitable acceptance rate.

G. CONVERGENCE DIAGNOSIS

It is important to ensure that the MCMC sampling is adequately exploring the parameter space and constructing an

accurate picture of the posterior distribution. One method of monitoring the performance of the adopted MCMC sampler is to examine *convergence diagnostics* that assess/monitor the extent to which the Markov chains have become a stationary distribution. Practitioners routinely apply the Gelman-Rubin (GR) convergence diagnostic [128] that is developed by sampling from multiple MCMC chains, whereby the variance of each chain is assessed independently (within-chain variance) and then compared to the variance between the multiple chains (between-chain variance) for each parameter. A large difference between these two variances would indicate that the chains have not converged on the same stationary distribution.

In our case, we run several independent experiments and compare the MCMC chains using a modified Gelman-Rubin convergence diagnostic presented by Vehtari et al. [129]. It is beyond the scope of this publication to provide a detailed mathematical description of the convergence diagnostics, the reader can refer to [129] for a full description. A useful package for Bayesian model diagnostics, which contains an implementation of this modified diagnostic is *arviz* [130]¹⁴. We refer to the modified Gelman-Rubin diagnostic as \hat{R} , where values close to 1 indicate convergence. In Listing 15, we present the code to prepare the MCMC sampler outputs

¹⁴<https://python.arviz.org/en/stable/api/generated/arviz.rhat.html>

```

1  n_accept = 0
2  n_langevin = 0
3  # Run the MCMC sample for n_samples
4  for ii in tqdm(np.arange(1,self.n_samples)):
5      # Sample new values for theta
6      theta_proposal = theta + np.random.normal(0, self.step_theta, self.theta_size)
7      lx = np.random.uniform(0,1,1)
8      if (self.use_langevin_gradients is True) and (lx < self.l_prob):
9          theta_gd = self.model.langevin_gradient(self.x_data, self.y_data, theta.copy(), self.
sgd_depth)
10         theta_proposal = np.random.normal(theta_gd, self.step_theta, self.theta_size)
11         theta_proposal_gd = self.model.langevin_gradient(self.x_data, self.y_data, theta_proposal.
copy(), self.sgd_depth)
12         # for numerical reasons, we will provide a simplified implementation that simplifies
13         # the MVN of the proposal distribution
14         wc_delta = (theta - theta_proposal_gd)
15         wp_delta = (theta_proposal - theta_gd)
16         sigma_sq = self.step_theta
17         first = -0.5 * np.sum(wc_delta * wc_delta) / sigma_sq
18         second = -0.5 * np.sum(wp_delta * wp_delta) / sigma_sq
19         diff_prop = first - second
20         n_langevin += 1
21     else:
22         diff_prop = 0
23         theta_proposal = np.random.normal(theta, self.step_theta, self.theta_size)
24         # calculate the prior
25         prior_proposal = self.prior(sigma_squared, theta_proposal) # takes care of the gradients
26         # calculate the likelihood considering observations
27         [likelihood_proposal, pred_y[ii,], sim_y[ii,], accuracy_data[ii]] = self.likelihood_function(
theta_proposal)
28         # calculate the test likelihood
29         [_, test_pred_y[ii,], test_sim_y[ii,], test_accuracy_data[ii]] = self.likelihood_function(
theta_proposal, test=True)
30         # since we using log scale: based on https://www.rapidtables.com/math/algebra/Logarithm.html
31         diff_likelihood = likelihood_proposal - likelihood
32         diff_prior = prior_proposal - prior_val
33         mh_prob = min(1, np.exp(diff_likelihood + diff_prior + diff_prop))
34         u = np.random.uniform(0, 1)
35         # Accept/reject
36         if u < mh_prob:
37             # Update position
38             n_accept += 1
39             # update
40             likelihood = likelihood_proposal
41             prior_val = prior_proposal
42             theta = theta_proposal
43             # and store
44             pos_theta[ii,] = theta_proposal
45         else:
46             # store
47             pos_theta[ii,] = pos_theta[ii-1,]
48     # print the % of times the proposal was accepted
49     accept_ratio = (n_accept / self.n_samples) * 100
50     print('{:.3}% were accepted'.format(accept_ratio))
51     # store the posterior of theta
52     self.pos_theta = pos_theta[self.n_burnin:, ]
53     # Create a pandas dataframe to store the posterior samples of theta
54     results_dict = {'w{}'.format(_): self.pos_theta[:, _].squeeze() for _ in range(self.theta_size-2)}
55     results_dict['b0'] = self.pos_theta[:, self.theta_size-2].squeeze()
56     results_dict['b1'] = self.pos_theta[:, self.theta_size-1].squeeze()
57     results_df = pd.DataFrame.from_dict(results_dict)
58     return results_df
59

```

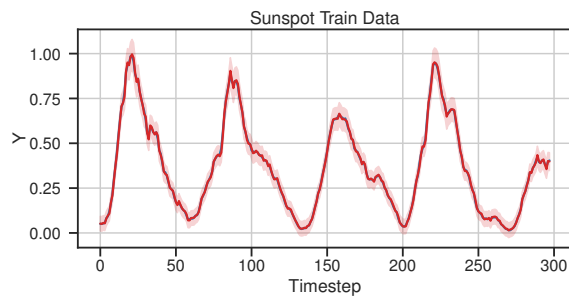
Listing 14: Begin with sampling loop (Continued from Listing 13)

Method	Problem	Train RMSE	Test RMSE	Accept. rate
Bayesian linear model	Sunspot	0.025 (0.013)	0.022 (0.012)	13.5%
Bayesian neural network	Sunspot	0.027 (0.007)	0.026 (0.007)	7.4%
Bayesian linear model	Abalone	0.085 (0.005)	0.086 (0.005)	5.8%
Bayesian neural network	Abalone	0.080 (0.002)	0.080 (0.002)	3.8%

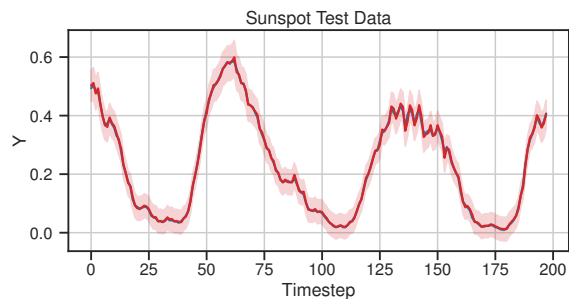
TABLE 1: Results using Bayesian linear model and Bayesian neural networks via MCMC sampler for the Abalone (regression) and the Sunspot (prediction) problems. The results show the RMSE mean and standard deviation (in brackets) for the train and test datasets, respectively.

Method	Problem	Train Accuracy	Test Accuracy	Accept. rate
Bayesian linear model	Iris	90.392% (2.832)	90.844% (3.039)	83.5%
Bayesian neural network	Iris	97.377% (0.655)	98.116% (1.657)	97.0%
Bayesian linear model	Ionosphere	89.060% (1.335)	85.316% (2.390)	58.8%
Bayesian neural network	Ionosphere	99.632% (0.356)	92.668% (1.890)	94.5%

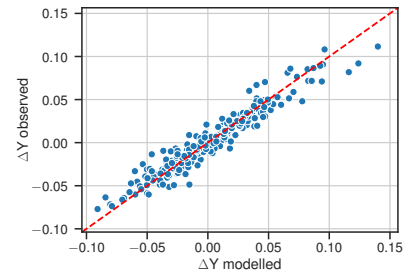
TABLE 2: Classification accuracy with Bayesian linear model and Bayesian neural networks via MCMC. The results show the accuracy mean and standard deviation (in brackets) for the train and test datasets.



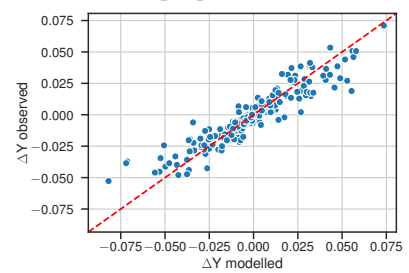
(a) Predictions and actual train dataset for Sunspot



(c) Predictions and actual test dataset for Sunspot

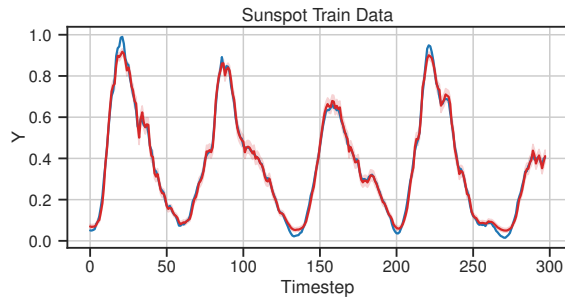


(b) Predictions and actual of ΔY train dataset for Sunspot prediction.

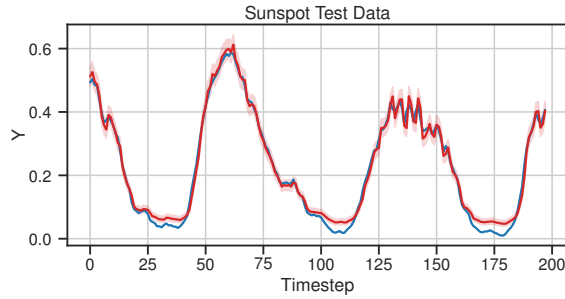


(d) Predictions and actual of ΔY test dataset for Sunspot prediction.

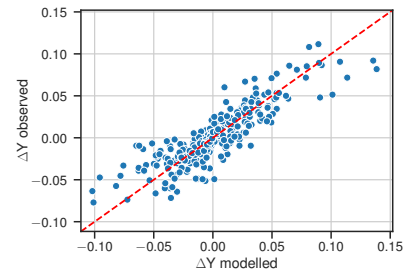
FIGURE 8: Bayesian linear regression model - Sunspot dataset



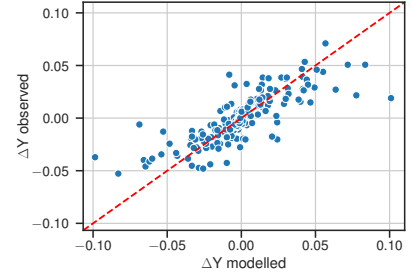
(a) Predictions and actual train dataset for Sunspot



(c) Predictions and actual test dataset for Sunspot



(b) Predictions and actual of ΔY train dataset for Sunspot prediction.



(d) Predictions and actual of ΔY test dataset for Sunspot prediction.

FIGURE 9: Bayesian neural network model - Sunspot dataset

for convergence diagnostic by *arviz*.

1) Results

We show results for the modified Gelman-Rubin diagnostic for Bayesian linear and Bayesian neural network models for each of the four datasets. We test five chains for each model, each with 5,000 samples of weights (excluding 50% burn-in samples) and computing the \hat{R} values for each parameter. We also provide an additional example "Linear+" for the linear models where each of the 5 chains has 50,000 samples (excluding 50% burn-in samples). The basic sampler presented here is not state of the art in terms of sampling efficiency (see Section -I for further discussion), and for some of the problems (particularly those where parameters are difficult to identify) it may take a large number of samples to converge. This additional example is shown to demonstrate that the convergence is improving as the number of samples grows. Figure 10 shows the distribution of the \hat{R} values, and we observe that the \hat{R} values of the weights for the Bayesian linear regression model are much smaller than Bayesian neural network. We can observe that based on the Gelman-Rubin diagnostics, the Bayesian neural network shows poor convergence. The additional samples in the "Linear+" case improve convergence, in particular for the Abalone, Iris and Ionosphere cases.

By closely examining each weight individually, we observe that the problem of non-convergence mainly arises from the multi-modality of the posterior. In Figure 11, we look at samples from a single chain of 20,000 samples

excluding 20% burn-in. In Figure 11a and 11b, we present a visualization for a selected weight from the Bayesian linear model. In Figure 11c and 11d, we present a visualization for a selected weight from the Bayesian neural network model. We observe potential multi-modal distributions in both cases, with a high degree of auto-correlation and poor convergence. To examine the impact of longer MCMC chains in achieving convergence, we ran an additional test by taking 400,000 samples excluding 20% burn-in, and then thinning the chain by a factor of 50 for visualisation. Thinning is the process of reducing the memory burden of the chain, particularly where samples may be auto-correlated [131]. We can also use thinning to more easily visualise long chains with many samples. In our case, we implement thinning by retaining every 50th sample of the chain. We present these results for the Iris dataset in Figure 12. We can see that the chains exhibit more desirable properties, particularly in the case of the linear model (Figure 12- Panels a and b). We can see in Figure 12 - Panels c and d, that the Bayesian neural network exhibits a multimodal posterior for this parameter.

Furthermore, other approaches for convergence diagnosis such as autocorrelation analysis can also be implemented with packages such as the *integrated autocorrelation time* [132] in *Emcee* [133]¹⁵. We refer to readers for a comprehensive review of MCMC convergence diagnostics given by [134]–[136].

¹⁵<https://emcee.readthedocs.io/en/stable/tutorials/autocorr/>

```

1 import arviz as az
2 # Sample a second chain as more than one chain is required to generate the Rhat diagnostic
3 # Run second chain using a different seed
4 np.random.seed(2)
5 # setup second chain
6 mcmc_chain2 = MCMC(n_samples, burn_in, x_data, y_data, x_test, y_test)
7 # Run the sampler
8 results_chain2 = mcmc_chain2.sampler()
9
10 # Now combine these two chains by stacking and convert to arviz dataset
11 # To convert to arviz, convert the pandas dataframe of samples into a dictionary of lists, stack and
12 # ingest into arviz
13 res_dict = results.to_dict(orient='list')
14 res_dict_chain2 = results_chain2.to_dict(orient='list')
15 # stack the chains looping through each parameter and ingest into arviz
16 az_results = az.from_dict(
17     {par: np.vstack([res_dict[par], res_dict_chain2[par]]) for par in res_dict}
18 )
19
20 # now we can use arviz to obtain summary statistics with Rhat ("r_hat") for each parameter.
21 az.summary(az_results)

```

Listing 15: Convergence diagnostics using the *arviz* python package that could be run after Listing 7

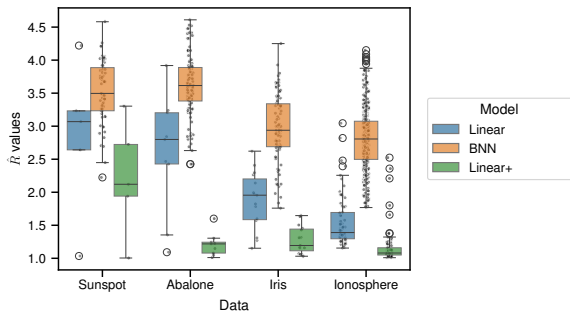


FIGURE 10: Distribution of \hat{R} values for our Bayesian linear model (Linear), Bayesian neural network (BNN) and a linear model with 10 times more samples (Linear+). The scattered points show the underlying data summarised in the box plots. The "Linear+" case shows the ability of the sampler to converge as more samples are taken.

H. MCMC PACKAGES

We note that there are avenues to further improve the sampling efficiency, we utilized Langevin MCMC but other gradient-based approaches such as Hamiltonian MCMC (HMC) [60] also exist. It is worthwhile to evaluate the performance of Langevin MCMC against other advanced gradient-based sampling algorithms (e.g., No-U-turn Sampler (NUTS) [137]) to assess convergence properties in cases of multimodal posterior distributions, such as in Bayesian neural networks. We note that implementation of HMC and NUTS exist in probabilistic programming libraries such as PyMC [138] and Stan [139], [140]. We implemented an additional notebook using the NumPyro [141] probabilistic programming library to perform an equivalent Bayesian lin-

ear regression (Listing 3-7) as an example ¹⁶.

I. DISCUSSION

We presented a Python tutorial for Bayesian neural networks and Bayesian linear models using MCMC sampling. In general, we observed that the Bayesian neural network performs better than Bayesian linear regression in our selected problems (Tables 1 and 2), despite showing no poor convergence (Figure 10, Figures and 12-Panel d). This could be due to the challenge of sampling a relatively large number of weights and biases of Bayesian neural networks, which also have multimodal posterior distribution. Hence, we conclude that for the case of Bayesian neural networks, a poor performance in the Gelman-Rubin diagnostics does not necessarily imply a poor performance in prediction tasks. We revisit the principle of *equifinality* [142], [143] which states that in open systems, a given end state can be reached by many potential means. In our case, the system is a neural network model and many solutions exist that represent a trained model displaying an accepted level of performance accuracy. However, we note that Bayesian models offer uncertainty quantification in predictions, and proper convergence is required. The original Gelman-Rubin diagnosis [128] motivated several enhancements for different types of problems [129], [135], [144], [145]; and we may need to develop a better diagnosis for Bayesian neural networks. Nonetheless, in our case comparing Bayesian logistic regression (converged) and Bayesian neural networks (not converged but achieved good accuracy), we can safely state that the Bayesian neural network presented can only provide a means for uncertainty quantification, but is not mature enough to qualify as a robust Bayesian model.

¹⁶https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial/blob/main/05-Linear-Model_NumPyro.ipynb

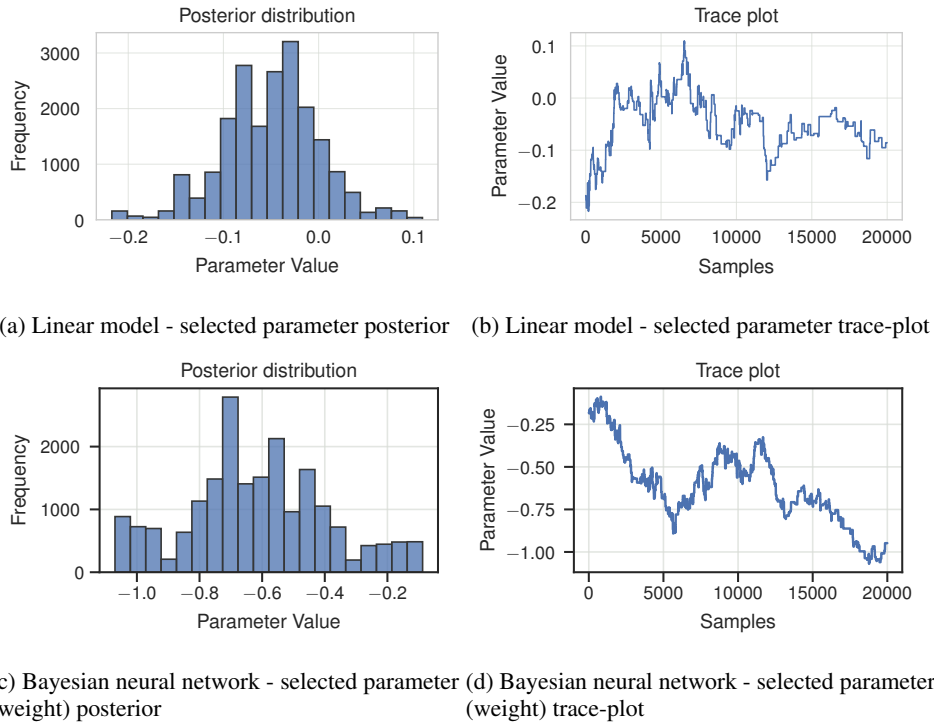


FIGURE 11: Posterior and trace-plot for a parameter in each of the Bayesian linear and Bayesian neural network models - Sunspot data. In this example, 20,000 samples were taken excluding 20% burn-in with no thinning.

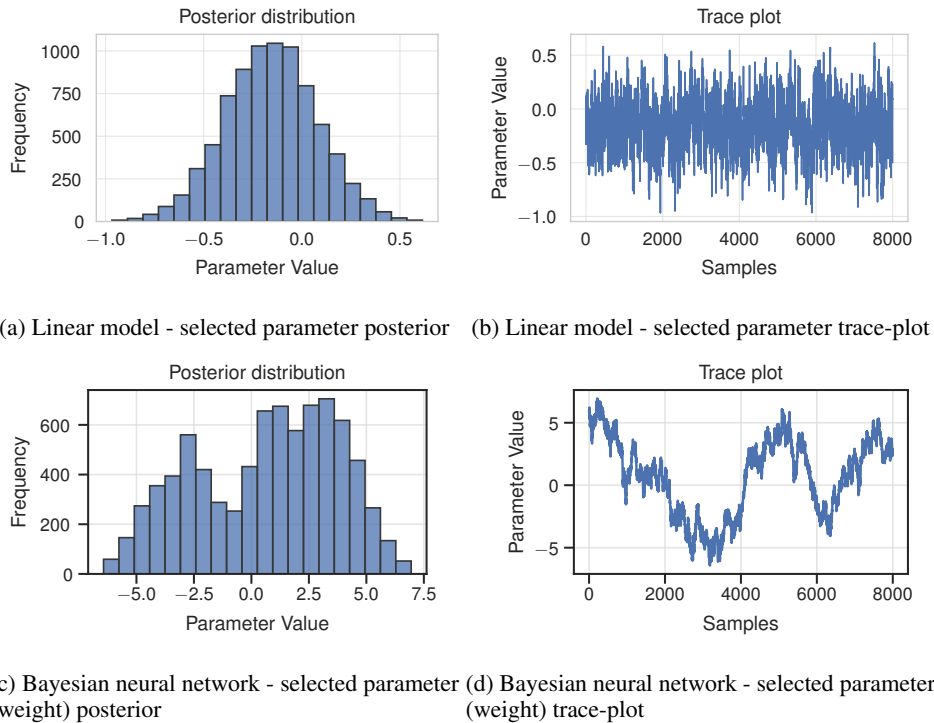


FIGURE 12: Posterior and trace-plot for a selected parameter in each of the Bayesian linear and Bayesian neural network models - Iris data. We took 400,000 samples (excluding 20% burn-in) with the chain thinned by a factor of 50 for visualisation.

We also revisit the convergence issue in the case of the Bayesian linear model as shown in Panel (b) - Figure 11, which shows a multi-modal distribution that has not well converged. It may also be the case that certain features are not contributing much to the decision-making process (predictions) and the weights associated (coefficients) with those features may be difficult to sample. This is similar to the case of neural networks, where certain weight links are not needed and can be pruned. Pruning neural networks create compact models that also get better generalization performance [146].

We note that we attained a much higher acceptance rate in Bayesian neural networks for classification problems (Table 2) when compared to regression problems (Table 2). We note that different likelihood functions are used for classification and regression (Multinomial and Gaussian likelihood) and we utilized the Langevin-gradient proposal distributions that accounted for the detailed balanced condition. Hence, we need to further fine-tune the hyperparameters associated with the proposal distribution to ensure we get a higher acceptance rate for the regression problems. Fine-tuning the hyperparameters for the proposal distribution is a laborious task which can be seen as a major limitation of MCMC sampling in large models such as Bayesian neural networks. Although 23 % acceptance rate [126] has been prominently used as a "golden rule", the optimal acceptance rate depends on the nature of the problem. The number of parameters, Langevin-based proposal distribution, and type of model would raise questions about the established acceptance rate [147]. Hence, more work needs to be done to establish what acceptance rates are appropriate for simple neural networks and deep learning models.

A way to address the issue of convergence would be to develop an ensemble of linear models that can compete with the accuracy of neural networks or deep learning models. In ensemble methods such as bagging and boosting, we can use linear models that have attained convergence as per Gelman-Rubin diagnosis and then combine the results of the ensemble using averaging and voting, as done in ensemble methods.

Our previous work has shown that despite the challenges, the combination of Langevin-gradients with parallel tempering MCMC [62], presents opportunities for sampling larger neural network architectures such as autoencoders and graph-based CNNs [69], [70]. The need to feature a robust methodology for uncertainty quantification in CNNs will make them more suitable for applications where uncertainty in decision-making poses major risks, such as medical image analysis [148] and human security [149]. CNNs have been considered for modelling temporal sequences, they have proven to be successful for time series classification [150], [151], and time series forecasting problems [152]–[154]. It has also been shown that one-dimensional CNNs provide better prediction performance than the conventional LSTM network for multi-step ahead time series prediction problems [154]. Leveraging CNNs within a Bayesian framework can provide better uncertainty quantification in predictions and make them useful for cutting-edge real-world applications.

We need a comprehensive evaluation of prominent gradient-based MCMC sampling methods for deep learning models such as CNNs, autoencoders, and LSTM networks.

We envision that this tutorial will enable statisticians, machine learning and deep learning experts to utilize MCMC sampling more effectively when developing new models and Bayesian frameworks for existing deep learning models. The tutorial has introduced basic concepts with code and provides an overview of challenges when it comes to the convergence of Bayesian neural networks. It can be further extended to utilize parallel computing via tempered MCMC [62], HMC for Bayesian neural networks, and Langevin MCMC and HMC for Bayesian deep learning.

CODE AND DATA

All code (for implementation, results and figures) and data presented in this paper are available in the associated github repository¹⁷. This repository presents the implementations in separate Jupyter notebooks in the base directory, with sub-directories containing data, convenient functions and details of the environment setup.

ACKNOWLEDGEMENTS AND CONTRIBUTIONS

R. Chandra contributed by conceptualisation, coding and experiments, analysis, writing, and project supervision. J. Simmons contributed by coding, analysis, and writing.

J. Simmons is funded by the Australian Research Council ITTC in Data Analytics for Resources and Environments (Grant IC19010031). We thank Royce Chen from UNSW Sydney for the initial experiments and analysis.

REFERENCES

- [1] G. Chamberlain and G. W. Imbens, "Nonparametric applications of bayesian inference," *Journal of Business & Economic Statistics*, vol. 21, no. 1, pp. 12–18, 2003.
- [2] J. Geweke, "Bayesian inference in econometric models using monte carlo integration," *Econometrica: Journal of the Econometric Society*, pp. 1317–1339, 1989.
- [3] G. E. Box and G. C. Tiao, *Bayesian inference in statistical analysis*, vol. 40. John Wiley & Sons, 2011.
- [4] R. Chandra, R. D. Müller, D. Azam, R. Deo, N. Butterworth, T. Salles, and S. Cripps, "Multicore parallel tempering Bayeslands for basin and landscape evolution," *Geochemistry, Geophysics, Geosystems*, vol. 20, no. 11, pp. 5082–5104, 2019.
- [5] J. Pall, R. Chandra, D. Azam, T. Salles, J. M. Webster, R. Scalzo, and S. Cripps, "Bayesreef: A bayesian inference framework for modelling reef growth in response to environmental change and biological dynamics," *Environmental Modelling & Software*, p. 104610, 2020.
- [6] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [7] L. Martino and V. Elvira, "Metropolis sampling," *Wiley StatsRef: Statistics Reference Online*, pp. 1–18, 2014.
- [8] D. J. MacKay, "Probable networks and plausible predictions—a review of practical bayesian methods for supervised neural networks," *Network: Computation in Neural Systems*, vol. 6, no. 3, pp. 469–505, 1995.
- [9] R. M. Neal, *Bayesian learning for neural networks*, vol. 118. Springer Science & Business Media, 2012.
- [10] D. F. Specht, "Probabilistic neural networks," *Neural networks*, vol. 3, no. 1, pp. 109–118, 1990.

¹⁷<https://github.com/sydney-machine-learning/Bayesianneuralnetworks-MCMC-tutorial>

- [11] M. D. Richard and R. P. Lippmann, "Neural network classifiers estimate bayesian a posteriori probabilities," *Neural computation*, vol. 3, no. 4, pp. 461–483, 1991.
- [12] E. A. Wan, "Neural network classification: A bayesian interpretation," *IEEE Transactions on Neural Networks*, vol. 1, no. 4, pp. 303–305, 1990.
- [13] Z. Ghahramani, "Probabilistic machine learning and artificial intelligence," *Nature*, vol. 521, no. 7553, pp. 452–459, 2015.
- [14] S. Chib and E. Greenberg, "Understanding the Metropolis-Hastings algorithm," *The american statistician*, vol. 49, no. 4, pp. 327–335, 1995.
- [15] C. P. Robert and G. Casella, "The Metropolis—Hastings algorithm," in *Monte Carlo statistical methods*, pp. 231–283, Springer, 1999.
- [16] D. B. Hitchcock, "A history of the Metropolis–Hastings algorithm," *The American Statistician*, vol. 57, no. 4, pp. 254–257, 2003.
- [17] G. Casella and E. I. George, "Explaining the Gibbs sampler," *The American Statistician*, vol. 46, no. 3, pp. 167–174, 1992.
- [18] C. K. Carter and R. Kohn, "On Gibbs sampling for state space models," *Biometrika*, vol. 81, no. 3, pp. 541–553, 1994.
- [19] G. O. Roberts and A. F. Smith, "Simple conditions for the convergence of the Gibbs sampler and Metropolis-Hastings algorithms," *Stochastic processes and their applications*, vol. 49, no. 2, pp. 207–216, 1994.
- [20] P. J. Rossky, J. D. Doll, and H. L. Friedman, "Brownian dynamics as smart Monte Carlo simulation," *The Journal of Chemical Physics*, vol. 69, no. 10, pp. 4628–4633, 1978.
- [21] G. O. Roberts and R. L. Tweedie, "Exponential convergence of langevin distributions and their discrete approximations," *Bernoulli*, pp. 341–363, 1996.
- [22] G. O. Roberts and J. S. Rosenthal, "Optimal scaling of discrete approximations to langevin diffusions," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 60, no. 1, pp. 255–268, 1998.
- [23] B. D. Flury, "Acceptance–rejection sampling made easy," *Siam Review*, vol. 32, no. 3, pp. 474–476, 1990.
- [24] W. R. Gilks and P. Wild, "Adaptive rejection sampling for Gibbs sampling," *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, vol. 41, no. 2, pp. 337–348, 1992.
- [25] S. T. Tokdar and R. E. Kass, "Importance sampling: a review," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, no. 1, pp. 54–60, 2010.
- [26] F. Llorente, L. Martino, D. Delgado-Gómez, and G. Camps-Valls, "Deep importance sampling based on regression for model inversion and emulation," *Digital Signal Processing*, vol. 116, p. 103104, 2021.
- [27] A. Brockwell, P. Del Moral, and A. Doucet, "Sequentially interacting Markov chain Monte Carlo methods," *The Annals of Statistics*, pp. 3387–3411, 2010.
- [28] C. Andrieu and J. Thoms, "A tutorial on adaptive MCMC," *Statistics and computing*, vol. 18, no. 4, pp. 343–373, 2008.
- [29] R. H. Swendsen and J.-S. Wang, "Replica monte carlo simulation of spin-glasses," *Physical review letters*, vol. 57, no. 21, p. 2607, 1986.
- [30] K. Hukushima and K. Nemoto, "Exchange monte carlo method and application to spin glass simulations," *Journal of the Physical Society of Japan*, vol. 65, no. 6, pp. 1604–1608, 1996.
- [31] D. J. Earl and M. W. Deem, "Parallel tempering: Theory, applications, and new perspectives," *Physical Chemistry Chemical Physics*, vol. 7, no. 23, pp. 3910–3916, 2005.
- [32] M. Sambridge, "A parallel tempering algorithm for probabilistic sampling and multimodal optimization," *Geophysical Journal International*, vol. 196, no. 1, pp. 357–374, 2014.
- [33] Y. Fan and S. A. Sisson, "Reversible jump MCMC," *Handbook of Markov Chain Monte Carlo*, pp. 67–92, 2011.
- [34] S. P. Brooks, P. Giudici, and G. O. Roberts, "Efficient construction of reversible jump Markov chain Monte Carlo proposal distributions," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 65, no. 1, pp. 3–39, 2003.
- [35] C. Tarantola, "Mcmc model determination for discrete graphical models," *Statistical Modelling*, vol. 4, no. 1, pp. 39–61, 2004.
- [36] F. Liang and C. Liu, "Efficient MCMC estimation of discrete distributions," *Computational statistics & data analysis*, vol. 49, no. 4, pp. 1039–1052, 2005.
- [37] G. Zanella, "Informed proposals for local mcmc in discrete spaces," *Journal of the American Statistical Association*, vol. 115, no. 530, pp. 852–865, 2020.
- [38] W. J. Browne, "MCMC algorithms for constrained variance matrices," *Computational Statistics & Data Analysis*, vol. 50, no. 7, pp. 1655–1677, 2006.
- [39] A. R. Gallant, H. Hong, M. P. Leung, and J. Li, "Constrained estimation using penalization and MCMC," *Journal of Econometrics*, vol. 228, no. 1, pp. 85–106, 2022.
- [40] S. A. Sisson and Y. Fan, *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839], 2011.
- [41] D. A. Van Dyk and X.-L. Meng, "The art of data augmentation," *Journal of Computational and Graphical Statistics*, vol. 10, no. 1, pp. 1–50, 2001.
- [42] L. L. Duan, J. E. Johndrow, and D. B. Dunson, "Scaling up data augmentation MCMC via calibration," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2575–2608, 2018.
- [43] J. Zobitz, A. Desai, D. Moore, and M. Chadwick, "A primer for data assimilation with ecological models using Markov Chain Monte Carlo (MCMC)," *Oecologia*, vol. 167, no. 3, pp. 599–611, 2011.
- [44] C. Andrieu, P. Djurić, and A. Doucet, "Model selection by MCMC computation," *Signal Processing*, vol. 81, no. 1, pp. 19–37, 2001.
- [45] C. Andrieu and A. Doucet, "Joint Bayesian model selection and estimation of noisy sinusoids via reversible jump MCMC," *IEEE Transactions on Signal Processing*, vol. 47, no. 10, pp. 2667–2676, 1999.
- [46] A. S. Mugglin, B. P. Carlin, L. Zhu, and E. Conlon, "Bayesian areal interpolation, estimation, and smoothing: an inferential approach for geographic information systems," *Environment and Planning A*, vol. 31, no. 8, pp. 1337–1352, 1999.
- [47] M. Sambridge, "Geophysical inversion with a neighbourhood algorithm—ii. appraising the ensemble," *Geophysical Journal International*, vol. 138, no. 3, pp. 727–746, 1999.
- [48] M. Sambridge and K. Mosegaard, "Monte carlo methods in geophysical inverse problems," *Reviews of Geophysics*, vol. 40, no. 3, 2002.
- [49] R. Scalzo, D. Kohn, H. Olierook, G. Houseman, R. Chandra, M. Girolami, and S. Cripps, "Efficiency and robustness in Monte Carlo sampling for 3-D geophysical inversions with Obsidian v0. 1.2: Setting up for success," *Geoscientific Model Development*, vol. 12, no. 7, pp. 2941–2960, 2019.
- [50] R. Chandra, D. Azam, R. D. Müller, T. Salles, and S. Cripps, "Bayeslands: A Bayesian inference approach for parameter uncertainty quantification in badlands," *Computers & Geosciences*, vol. 131, pp. 89–101, 2019.
- [51] H. K. Olierook, R. Scalzo, D. Kohn, R. Chandra, E. Farahbakhsh, C. Clark, S. M. Reddy, and R. D. Müller, "Bayesian geological and geophysical data fusion for the construction and uncertainty quantification of 3D geological models," *Geoscience Frontiers*, vol. 12, no. 1, pp. 479–493, 2021.
- [52] L. Marshall, D. Nott, and A. Sharma, "A comparative study of Markov chain Monte Carlo methods for conceptual rainfall-runoff modeling," *Water Resources Research*, vol. 40, no. 2, 2004.
- [53] J. A. Vrugt, C. J. Ter Braak, M. P. Clark, J. M. Hyman, and B. A. Robinson, "Treatment of input uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo simulation," *Water Resources Research*, vol. 44, no. 12, 2008.
- [54] G. I. Valderrama-Bahamóndez and H. Fröhlich, "MCMC techniques for parameter estimation of ODE based models in systems biology," *Frontiers in Applied Mathematics and Statistics*, vol. 5, p. 55, 2019.
- [55] Y. Nishiyama, Y. Saikawa, and N. Nishiyama, "Interaction between the immune system and acute myeloid leukemia: A model incorporating promotion of regulatory T cell expansion by leukemic cells," *Biosystems*, vol. 165, pp. 99–105, 2018.
- [56] B. Rannala, "Identifiability of parameters in MCMC Bayesian inference of phylogeny," *Systematic biology*, vol. 51, no. 5, pp. 754–760, 2002.
- [57] D. Sorensen, D. Gianola, and D. Gianola, "Likelihood, Bayesian and MCMC methods in quantitative genetics," 2002.
- [58] A. J. Drummond, M. A. Suchard, D. Xie, and A. Rambaut, "Bayesian phylogenetics with beautil and the beast 1.7," *Molecular biology and evolution*, vol. 29, no. 8, pp. 1969–1973, 2012.
- [59] M. Girolami and B. Calderhead, "Riemann manifold langevin and hamiltonian monte carlo methods," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 2, pp. 123–214, 2011.
- [60] R. M. Neal et al., "Mcmc using hamiltonian dynamics," *Handbook of Markov Chain Monte Carlo*, vol. 2, no. 11, 2011.
- [61] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 681–688, 2011.
- [62] R. Chandra, K. Jain, R. V. Deo, and S. Cripps, "Langevin-gradient parallel tempering for bayesian neural learning," *Neurocomputing*, vol. 359, pp. 315 – 326, 2019.

- [63] M. M. Dragan and D. Thierens, "Evolutionary Markov chain Monte Carlo," in *International Conference on Artificial Evolution (Evolution Artificielle)*, pp. 63–76, Springer, 2004.
- [64] C. J. Ter Braak, "A Markov Chain Monte Carlo version of the genetic algorithm differential evolution: easy Bayesian computing for real parameter spaces," *Statistics and Computing*, vol. 16, no. 3, pp. 239–249, 2006.
- [65] C. J. Ter Braak and J. A. Vrugt, "Differential evolution markov chain with snooker updater and fewer chains," *Statistics and Computing*, vol. 18, no. 4, pp. 435–446, 2008.
- [66] A. Kapoor, E. Nukala, and R. Chandra, "Bayesian neuroevolution using distributed swarm optimization and tempered MCMC," *Applied Soft Computing*, vol. 129, p. 109528, 2022.
- [67] L. Li, A. Holbrook, B. Shahbaba, and P. Baldi, "Neural network gradient hamiltonian monte carlo," *Computational statistics*, vol. 34, no. 1, pp. 281–299, 2019.
- [68] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*, pp. 177–186, Springer, 2010.
- [69] R. Chandra, M. Jain, M. Maharana, and P. N. Krivitsky, "Revisiting Bayesian autoencoders with MCMC," *IEEE Access*, vol. 10, pp. 40482–40495, 2022.
- [70] R. Chandra, A. Bhagat, M. Maharana, and P. N. Krivitsky, "Bayesian graph convolutional neural networks via tempered MCMC," *IEEE Access*, vol. 9, pp. 130353–130365, 2021.
- [71] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, "Variational inference: A review for statisticians," *Journal of the American Statistical Association*, vol. 112, no. 518, pp. 859–877, 2017.
- [72] A. Graves, "Practical variational inference for neural networks," in *Advances in Neural Information Processing Systems 24*, pp. 2348–2356, 2011.
- [73] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural network," in *Proceedings of The 32nd International Conference on Machine Learning*, pp. 1613–1622, 2015.
- [74] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [75] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*, pp. 1050–1059, 2016.
- [76] K. Shridhar, F. Laumann, and M. Liwicki, "A comprehensive guide to bayesian convolutional neural network with variational inference," 2019.
- [77] Y. Gal and Z. Ghahramani, "A theoretically grounded application of dropout in recurrent neural networks," *Advances in neural information processing systems*, vol. 29, 2016.
- [78] J. M. Bernardo and A. F. Smith, *Bayesian theory*, vol. 405. John Wiley & Sons, 2009.
- [79] G. A. Barnard and T. Bayes, "Studies in the history of probability and statistics: IX. Thomas Bayes's essay towards solving a problem in the doctrine of chances," *Biometrika*, vol. 45, no. 3/4, pp. 293–315, 1958.
- [80] A. I. Dale, *A history of inverse probability: From Thomas Bayes to Karl Pearson*. Springer Science & Business Media, 2012.
- [81] P. S. Laplace, "Memoir on the probability of the causes of events," *Statistical science*, vol. 1, no. 3, pp. 364–378, 1986.
- [82] S. R. Johnson, G. A. Tomlinson, G. A. Hawker, J. T. Granton, and B. M. Feldman, "Methods to elicit beliefs for Bayesian priors: a systematic review," *Journal of clinical epidemiology*, vol. 63, no. 4, pp. 355–369, 2010.
- [83] K. M. Banner, K. M. Irvine, and T. J. Rodhouse, "The use of Bayesian priors in ecology: The good, the bad and the not great," *Methods in Ecology and Evolution*, vol. 11, no. 8, pp. 882–889, 2020.
- [84] R. van de Schoot, S. Depaoli, R. King, B. Kramer, K. Mårtens, M. G. Tadesse, M. Vannucci, A. Gelman, D. Veen, J. Willemsen, et al., "Bayesian statistics and modelling," *Nature Reviews Methods Primers*, vol. 1, no. 1, pp. 1–26, 2021.
- [85] W. Li and G. Lin, "An adaptive importance sampling algorithm for bayesian inversion with multimodal distributions," *Journal of Computational Physics*, vol. 294, pp. 173–190, 2015.
- [86] S. Hu, D. S. Poskitt, and X. Zhang, "Bayesian adaptive bandwidth kernel density estimation of irregular multivariate distributions," *Computational Statistics & Data Analysis*, vol. 56, no. 3, pp. 732–740, 2012.
- [87] R. Rojas, "A short proof of the posterior probability property of classifier neural networks," *Neural Computation*, vol. 8, no. 1, pp. 41–43, 1996.
- [88] H. Akaike, "Likelihood and the bayes procedure," in *Selected Papers of Hirotugu Akaike*, pp. 309–332, Springer, 1998.
- [89] D. J. MacKay, "Hyperparameters: Optimize, or integrate out?," in *Maximum entropy and Bayesian methods*, pp. 43–59, Springer, 1996.
- [90] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin, *Bayesian data analysis*. Boca Raton, Fla.: Chapman & Hall/CRC, 2004.
- [91] R. McElreath, *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Chapman and Hall/CRC, 2 ed., Mar. 2020.
- [92] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, "Array programming with NumPy," *Nature*, vol. 585, pp. 357–362, Sept. 2020.
- [93] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020.
- [94] G. O. Roberts and J. S. Rosenthal, "General state space Markov chains and MCMC algorithms," *Probability surveys*, vol. 1, pp. 20–71, 2004.
- [95] G. O. Roberts, J. S. Rosenthal, J. Segers, and B. Sousa, "Extremal indices, geometric ergodicity of Markov chains, and MCMC," *Extremes*, vol. 9, no. 3, pp. 213–229, 2006.
- [96] C. Andrieu and É. Moulines, "On the ergodicity properties of some adaptive MCMC algorithms," *The Annals of Applied Probability*, vol. 16, no. 3, pp. 1462–1505, 2006.
- [97] J. Grazzini, "Analysis of the emergent properties: Stationarity and ergodicity," *Journal of Artificial Societies and Social Simulation*, vol. 15, no. 2, p. 7, 2012.
- [98] B. Cheng and D. M. Titterton, "Neural networks: A review from a statistical perspective," *Statistical science*, pp. 2–30, 1994.
- [99] A. F. Smith, "A general Bayesian linear model," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 35, no. 1, pp. 67–75, 1973.
- [100] E. J. Bedrick, R. Christensen, and W. Johnson, "A new perspective on priors for generalized linear models," *Journal of the American Statistical Association*, vol. 91, no. 436, pp. 1450–1460, 1996.
- [101] V. Fortuin, "Priors in Bayesian deep learning: A review," *International Statistical Review*, 2022.
- [102] B. P. Hobbs, D. J. Sargent, and B. P. Carlin, "Commensurate priors for incorporating historical information in clinical trials using general and generalized linear models," *Bayesian Analysis*, vol. 7, no. 3, p. 639, 2012.
- [103] E. De Bézenac, A. Pajot, and P. Gallinari, "Deep learning for physical processes: Incorporating prior scientific knowledge," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2019, no. 12, p. 124009, 2019.
- [104] R. Chandra and Y. He, "Bayesian neural networks for stock price forecasting before and during COVID-19 pandemic," *PLOS ONE*, vol. 16, p. e0253217, July 2021.
- [105] A. Krogh and J. A. Hertz, "A simple weight decay can improve generalization," in *Advances in neural information processing systems*, pp. 950–957, 1992.
- [106] G. C. McDonald, "Ridge regression," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 1, no. 1, pp. 93–100, 2009.
- [107] S. I. Gallant et al., "Perceptron-based learning algorithms," *IEEE Transactions on neural networks*, vol. 1, no. 2, pp. 179–191, 1990.
- [108] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint arXiv:1609.04747*, 2016.
- [109] E. I. George, U. Makov, and A. Smith, "Conjugate likelihood distributions," *Scandinavian Journal of Statistics*, pp. 147–156, 1993.
- [110] K. P. Murphy, "Conjugate Bayesian analysis of the Gaussian distribution," *def*, vol. 1, no. 2, p. 16, 2007.
- [111] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [112] C. Bishop, C. M. Bishop, et al., *Neural networks for pattern recognition*. Oxford university press, 1995.
- [113] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.

- [114] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [115] B. J. Wythoff, "Backpropagation neural networks: a tutorial," *Chemometrics and Intelligent Laboratory Systems*, vol. 18, no. 2, pp. 115–155, 1993.
- [116] D. E. Rumelhart, R. Durbin, R. Golden, and Y. Chauvin, "Backpropagation: The basic theory," *Backpropagation: Theory, architectures and applications*, pp. 1–34, 1995.
- [117] A. Krogh and J. Hertz, "A simple weight decay can improve generalization," *Advances in neural information processing systems*, vol. 4, 1991.
- [118] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International conference on machine learning*, pp. 1139–1147, PMLR, 2013.
- [119] M. Moreira and E. Fiesler, "Neural networks with adaptive learning rate and momentum terms," *tech. rep., Idiap*, 1995.
- [120] R. Battiti, "First-and second-order methods for learning: between steepest descent and Newton's method," *Neural computation*, vol. 4, no. 2, pp. 141–166, 1992.
- [121] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for on-line learning and stochastic optimization," *Journal of Machine Learning Research*, vol. 12, no. Jul, pp. 2121–2159, 2011.
- [122] M. D. Zeiler, "Adadelta: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012.
- [123] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015.
- [124] H. Suwa and S. Todo, "Markov chain Monte Carlo method without detailed balance," *Physical review letters*, vol. 105, no. 12, p. 120603, 2010.
- [125] F. Takens, "Detecting strange attractors in turbulence," in *Dynamical Systems and Turbulence*, Warwick 1980, *Lecture Notes in Mathematics*, pp. 366–381, 1981.
- [126] A. Gelman, W. R. Gilks, and G. O. Roberts, "Weak convergence and optimal scaling of random walk Metropolis algorithms," *The annals of applied probability*, vol. 7, no. 1, pp. 110–120, 1997.
- [127] C. Andrieu and É. Moulines, "On the ergodicity properties of some adaptive MCMC algorithms," *The Annals of Applied Probability*, vol. 16, no. 3, pp. 1462 – 1505, 2006.
- [128] A. Gelman and D. B. Rubin, "Inference from iterative simulation using multiple sequences," *Statistical science*, pp. 457–472, 1992.
- [129] A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P.-C. Bürkner, "Rank-Normalization, Folding, and Localization: An Improved R[^] for Assessing Convergence of MCMC (with Discussion)," *Bayesian Analysis*, vol. 16, June 2021.
- [130] R. Kumar, C. Carroll, A. Hartikainen, and O. Martin, "ArviZ a unified library for exploratory analysis of Bayesian models in Python," *Journal of Open Source Software*, vol. 4, p. 1143, Jan. 2019.
- [131] W. A. Link and M. J. Eaton, "On thinning of chains in MCMC," *Methods in Ecology and Evolution*, vol. 3, pp. 112–115, Feb. 2012.
- [132] U. Wolff, A. Collaboration, et al., "Monte Carlo errors with less errors," *Computer Physics Communications*, vol. 156, no. 2, pp. 143–153, 2004.
- [133] D. Foreman-Mackey, D. W. Hogg, D. Lang, and J. Goodman, "emcee: the MCMC hammer," *Publications of the Astronomical Society of the Pacific*, vol. 125, no. 925, p. 306, 2013.
- [134] M. K. Cowles and B. P. Carlin, "Markov chain Monte Carlo convergence diagnostics: a comparative review," *Journal of the American statistical Association*, vol. 91, no. 434, pp. 883–904, 1996.
- [135] V. Roy, "Convergence diagnostics for Markov Chain Monte Carlo," *Annual Review of Statistics and Its Application*, vol. 7, pp. 387–412, 2020.
- [136] L. F. South, M. Riabiz, O. Teymur, and C. J. Oates, "Postprocessing of MCMC," *Annual Review of Statistics and Its Application*, vol. 9, pp. 529–555, 2022.
- [137] M. D. Hoffman and A. Gelman, "The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1593–1623, 2014.
- [138] A. Patil, D. Huard, and C. J. Fonnesbeck, "PyMC: Bayesian stochastic modelling in Python," *Journal of statistical software*, vol. 35, no. 4, p. 1, 2010.
- [139] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. A. Brubaker, J. Guo, P. Li, and A. Riddell, "Stan: A probabilistic programming language," *Journal of statistical software*, vol. 76, 2017.
- [140] J. Annis, B. J. Miller, and T. J. Palmeri, "Bayesian inference with Stan: a tutorial on adding custom distributions," *Behavior research methods*, vol. 49, pp. 863–886, 2017.
- [141] D. Phan, N. Pradhan, and M. Jankowiak, "Composable Effects for Flexible and Accelerated Probabilistic Programming in NumPyro," *Dec. 2019*.
- [142] D. Cicchetti and F. A. Rogosch, "Equifinality and multifinality in developmental psychopathology," *Development and psychopathology*, vol. 8, no. 4, pp. 597–600, 1996.
- [143] C. Gresov and R. Drazin, "Equifinality: Functional equivalence in organization design," *Academy of management review*, vol. 22, no. 2, pp. 403–428, 1997.
- [144] S. P. Brooks and A. Gelman, "General methods for monitoring convergence of iterative simulations," *Journal of computational and graphical statistics*, vol. 7, no. 4, pp. 434–455, 1998.
- [145] D. Vats and C. Knudson, "Revisiting the Gelman–Rubin diagnostic," *Statistical Science*, vol. 36, no. 4, pp. 518–529, 2021.
- [146] T. Liang, J. Glossner, L. Wang, S. Shi, and X. Zhang, "Pruning and quantization for deep neural network acceleration: A survey," *Neurocomputing*, vol. 461, pp. 370–403, 2021.
- [147] M. Bedard, "Optimal acceptance rates for Metropolis algorithms: Moving beyond 0.234," *Stochastic Processes and their Applications*, vol. 118, no. 12, pp. 2198–2222, 2008.
- [148] S. Anwar, M. Majid, and A. Qayyum, "Medical image analysis using convolutional neural networks: A review," *Journal of Medical Systems*, vol. 42, no. 226, 2018.
- [149] S. Akcay, M. E. Kundegorski, C. G. Willcocks, and T. P. Breckon, "Using deep convolutional neural network architectures for object classification and detection within x-ray baggage security imagery," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 9, pp. 2203–2215, 2018.
- [150] B. Zhao, H. Lu, S. Chen, J. Liu, and D. Wu, "Convolutional neural networks for time series classification," *Journal of Systems Engineering and Electronics*, vol. 28, no. 1, pp. 162–169, 2017.
- [151] C.-L. Liu, W.-H. Hsaio, and Y.-C. Tu, "Time series classification with multivariate convolutional neural network," *IEEE Transactions on Industrial Electronics*, vol. 66, no. 6, pp. 4788–4797, 2018.
- [152] A. Borovykh, S. Bohte, and C. W. Oosterlee, "Conditional time series forecasting with convolutional neural networks," *arXiv preprint arXiv:1703.04691*, 2017.
- [153] M. Binkowski, G. Marti, and P. Donnat, "Autoregressive convolutional neural networks for asynchronous time series," in *International Conference on Machine Learning*, pp. 580–589, PMLR, 2018.
- [154] R. Chandra, S. Goyal, and R. Gupta, "Evaluation of deep learning models for multi-step ahead time series prediction," *IEEE Access*, vol. 9, pp. 83105–83123, 2021.



DR ROHITASH CHANDRA is a Senior Lecturer in Data Science at the UNSW School of Mathematics and Statistics. He leads a program of research encircling methodologies and applications of artificial intelligence. The methodologies include Bayesian deep learning, neuroevolution, ensemble learning, and data augmentation. The applications include climate extremes, geoscientific models, mineral exploration, biomedicine, and language modelling. Dr Chandra is on the Editorial Board for Geoscientific Model Development. He also served as an Associate Editor for IEEE Transactions on Neural Networks and Learning Systems, and Neurocomputing (2021–2022). Since 2021, Dr Chandra has been featured in the Stanford's list of top 2% scientists in the world.



DR JOSHUA SIMMONS received his Ph.D from UNSW Sydney and is currently working as a Research Engineer at the ARC Training Centre in Data Analytics for Resources and Environments (DARE). His research has included numerical and statistical modelling of coastal processes, with a particular focus on storm erosion of sandy beaches. More recently, his work has focused on interpretable machine learning/data science approaches and Bayesian methods for uncertainty

quantification in the hydrology domain.

• • •