

A Flow-Based Generative Model for Rare-Event Simulation

Lachlan Gibson¹, Marcus Hoerger^{1*} and Dirk Kroese¹

^{1*}School of Mathematics & Physics, The University of
Queensland, Brisbane, 4072, QLD, Australia.

*Corresponding author(s). E-mail(s): m.hoerger@uq.edu.au;
Contributing authors: l.gibson1@uq.edu.au;
kroese@maths.uq.edu.au;

Abstract

Solving decision problems in complex, stochastic environments is often achieved by estimating the expected outcome of decisions via Monte Carlo sampling. However, sampling may overlook rare, but important events, which can severely impact the decision making process. We present a method in which a Normalizing Flow generative model is trained to simulate samples directly from a conditional distribution given that a rare event occurs. By utilizing Coupling Flows, our model can, in principle, approximate any sampling distribution arbitrarily well. By combining the approximation method with Importance Sampling, highly accurate estimates of complicated integrals and expectations can be obtained. We include several examples to demonstrate how the method can be used for efficient sampling and estimation, even in high-dimensional and rare-event settings. We illustrate that by simulating directly from a rare-event distribution significant insight can be gained into the way rare events happen.

Keywords: normalizing flows, neural networks, rare events, simulation

1 Introduction

Many decision problems in complex, stochastic environments ([Kochenderfer, 2015](#)) are nowadays solved by estimating the expected outcome of decisions via

Monte Carlo simulations (Kroese, Botev, Taimre, & Vaisman, 2019; Kroese, Taimre, & Botev, 2011; Liu, 2004). The success of Monte Carlo methods is due to their simplicity, flexibility, and scalability. However, in many problem domains, the occurrence of rare but important events — events that happen with a very small probability, say less than 10^{-4} — severely impairs their efficiency, for the reason that such events do not show up often in a typical simulation run (Arief et al., 2021). On the other hand, failing to consider such rare events can lead to decisions with potentially catastrophic outcomes, e.g., hazardous behaviour of autonomous cars. By using well-known variance reduction techniques such as Importance Sampling (Bucklew, 2004; Kroese et al., 2011) it is possible to sometimes dramatically increase the efficiency of the standard Monte Carlo method. Nevertheless, there are few efficient methods available that give insights into *how* a system behaves under a rare event. An important research goal is thus to find methods that simulate a random process *conditionally* on the occurrence of a rare event. In this case, the target sampling distribution is the distribution of the original process conditioned on the rare event occurring. More broadly, for both estimation and sampling problems, the challenge is to identify a “good” sampling distribution that closely approximates a target distribution. Certain approximation methods, such as the Cross-Entropy method (de Boer, Kroese, Mannor, & Rubinstein, 2005; Rubinstein & Kroese, 2017), train a parametric model by minimizing the cross-entropy between a distribution family and the target distribution. However, typical parametric models are often not flexible enough to efficiently capture all the complexity of many interesting systems. Botev, Kroese, and Taimre (2007) introduced a Generalized Cross-entropy method, which approximates the target distribution in a non-parametric way. However, the quality of the results is determined by various hyper-parameters, which require expertise to tune.

While many standard methods can identify sampling distributions of sufficient quality to estimate rare-event probabilities and related quantities of interest, extremely good approximations to the target conditional distribution are required to be able to explore the behaviour of the simulated system under rare-event conditions. This can occasionally be achieved by strategically selecting very specific distribution families using problem-specific information, but this has been difficult to achieve with conventional simulation methods. Therefore, it is desirable to have a more general approach that can closely approximate any target distribution without relying on problem-specific knowledge. Gibson and Kroese (2022) recently introduced a framework for rare-event simulation using neural networks that aimed to achieve this. In that framework two Multilayer Perceptrons are trained simultaneously: the first being a generative model to represent the sampling distribution, and the second to approximate the probability distribution of the first. While the framework had some success in one-dimensional problems, the reliance on a second network and probability density estimation convoluted the training process and made scaling to higher-dimensional problems difficult. Ardizzone, Lüth, Kruse,

Rother, and Köthe (2019) and Falkner, Coretti, Romano, Geissler, and Delgado (2022) explore the use of normalizing flows for rare event sampling in the context of image generation and Boltzmann generators. The approach in Ardizzone et al. (2019) requires a conditioning network to be pre-trained and embedded into the normalizing flows architecture. Falkner et al. (2022) embeds a bias variable into the flow to learn distributions that are biased towards the region of interest. However, such bias variables can be difficult to construct for more complex problems. In contrast, our model is trained end-to-end and without having to introduce a conditioning variable.

We present a new framework, inspired by Gibson and Kroese (2022), in which a Normalizing Flows generative model is trained to learn the optimal sampling distribution and used to estimate quantities such as the rare-event probability. The highly expressive nature of some Normalizing Flows architectures, trained using standard deep learning training algorithms, massively expands the range of learnable distributions, while the invertibility of Normalizing Flows allows the exact generative probability density to be computed without the need of a second network. In Section 2 we present the background theory and a description of the training algorithm, and in Section 3 we present several examples of rare-event simulation using this method. The source-code of our method is available at <https://github.com/hoergems/rare-event-simulation-normalizing-flows>.

2 Theory

The theory of rare-event simulation (Bucklew, 2004; Juneja & Shahabuddin, 2006; Rubino & Tuffin, 2009), Importance Sampling (Glynn & Iglehart, 1989; Neal, 2001; Tokdar & Kass, 2010) and Normalizing Flows (Kobyzev, Prince, & Brubaker, 2021; Papamakarios, Nalisnick, Rezende, Mohamed, & Lakshminarayanan, 2021; Rezende & Mohamed, 2015) is well established. In this section we review the relevant background, formulate the foundation of our method, and present the training procedure.

2.1 Rare-Event Simulation and Importance Sampling

First let us introduce a mathematical basis for rare-event simulation and Importance Sampling. Consider a random object (e.g., random variable, random vector, or stochastic process), \mathbf{X} , taking values in some space, \mathcal{X} , with probability density function, $p(\mathbf{x})$, that represents the result of a simulation experiment. Running simulations corresponds to sampling \mathbf{X} . We consider rare events of the form $\{\mathbf{X} \in \mathcal{A}\}$, where $\mathcal{A} \subset \mathcal{X}$ is the set of simulation outcomes \mathbf{x} that satisfy a predicate $S(\mathbf{x}) \geq \gamma$ for some *performance function* $S : \mathcal{X} \rightarrow \mathbb{R}$ and *level parameter* $\gamma \in \mathbb{R}$. The probability c of the rare event $\{\mathbf{X} \in \mathcal{A}\}$ can thus be written as

$$c := \mathbb{P}(\mathbf{X} \in \mathcal{A}) = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}[\mathbb{1}_{\mathcal{A}}(\mathbf{X})],$$

4 *A Flow-Based Generative Model for Rare-Event Simulation*

where c is very small, but not 0. Here, $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$ is an indicator function that is 1 when $\mathbf{x} \in \mathcal{A}$ and 0 otherwise, and \mathbb{E} represents the expected value. The expected number of simulations to sample a single rare event is $1/c$, so simulating rare events by sampling \mathbf{X} directly quickly becomes computationally infeasible the rarer the event is. Therefore, other techniques, such as Importance Sampling become necessary to simulate and analyse rare events.

In particular, suppose that \mathbf{X} is sampled from a probability density function $p(\mathbf{x})$ and that we wish to estimate the quantity

$$\ell := \mathbb{E}[H(\mathbf{X})] = \int H(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x}$$

via “crude” Monte Carlo: Simulate $\mathbf{X}_1, \dots, \mathbf{X}_n$ independently from p , and estimate ℓ via the sample average $\sum_{i=1}^n H(\mathbf{X}_i)/n$. The idea of Importance Sampling is to change the probability distribution under which the simulation takes place. However, computing the expected value of a function of \mathbf{X} while using a different sampling density q , requires a correction factor of the likelihood ratio $p(\mathbf{x})/q(\mathbf{x})$,

$$\ell := \mathbb{E}_p [H(\mathbf{X})] = \mathbb{E}_q \left[H(\mathbf{X}) \frac{p(\mathbf{X})}{q(\mathbf{X})} \right], \quad (1)$$

where $\mathbb{E}_p = \mathbb{E}$ and \mathbb{E}_q represent the expected values under the two probability models. This relationship allows an unbiased estimate of ℓ to be computed by sampling \mathbf{X} from q instead of p , via

$$\widehat{\ell} := \frac{1}{n} \sum_{k=1}^n H(\mathbf{X}_k) \frac{p(\mathbf{X}_k)}{q(\mathbf{X}_k)},$$

where $\mathbf{X}_1, \dots, \mathbf{X}_n$ is an independent and identically distributed (iid) sample from q . Note that any choice for q is allowed, as long as $q(\mathbf{x}) \neq 0$ when $p(\mathbf{x}) \neq 0$. The optimal sampling distribution for estimating ℓ in this way is thus the distribution that minimizes the variance of the estimator $\widehat{\ell}$. When H is strictly positive (or strictly negative), choosing the probability density

$$q^*(\mathbf{x}) := \frac{p(\mathbf{x})H(\mathbf{x})}{\ell}, \quad (2)$$

yields in fact a zero-variance estimator, because then $\widehat{\ell} = \ell$. As a special case, an unbiased estimator of the rare-event probability $c = \mathbb{P}(\mathbf{X} \in \mathcal{A})$ can be computed as

$$\widehat{c} := \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\mathcal{A}}(\mathbf{X}_k) \frac{p(\mathbf{X}_k)}{q(\mathbf{X}_k)}, \quad (3)$$

and the optimal sampling density q^* is just the original density p truncated to the rare-event region \mathcal{A} :

$$q^*(\mathbf{x}) = p^{\mathcal{A}}(\mathbf{x}) := \frac{p(\mathbf{x})\mathbb{1}_{\mathcal{A}}(\mathbf{x})}{c}. \quad (4)$$

More generally, if we want to estimate the expectation of $H(\mathbf{X})$ conditional on $\mathbf{X} \in \mathcal{A}$, that is,

$$\ell^{\mathcal{A}} := \mathbb{E}_{p^{\mathcal{A}}}[H(\mathbf{X})] = \frac{1}{c} \mathbb{E}_q \left[H(\mathbf{X}) \frac{p(\mathbf{X})}{q(\mathbf{X})} \mathbb{1}_{\mathcal{A}}(\mathbf{X}) \right], \quad (5)$$

then $\ell^{\mathcal{A}}$ can be estimated via

$$\widehat{\ell}^{\mathcal{A}} := \frac{1}{n \widehat{c}} \sum_{k=1}^n H(\mathbf{X}_k) \frac{p(\mathbf{X}_k)}{q(\mathbf{X}_k)} \mathbb{1}_{\mathcal{A}}(\mathbf{X}_k). \quad (6)$$

In this case, the optimal Importance Sampling density is

$$q^*(\mathbf{x}) := \frac{p(\mathbf{x})H(\mathbf{x})\mathbb{1}_{\mathcal{A}}(\mathbf{x})}{\ell c}, \quad (7)$$

2.2 Normalizing Flows

If it is possible to approximate the target density (e.g., the optimal Importance Sampling density) closely, then we are able to compute certain quantities of interest (e.g., $c = \mathbb{P}(\mathbf{X} \in \mathcal{A})$ or $\ell = \mathbb{E}H(\mathbf{X})$) with low variance. *Normalizing Flows* is a generative model that can closely approximate a wide variety of probability densities. Suppose the random object, \mathbf{X} , can be mapped to another random object, \mathbf{Z} , taking values in some space \mathcal{Z} of the same dimensionality, and vice versa, via invertible differentiable functions $\boldsymbol{\psi}$ and $\boldsymbol{\phi} = \boldsymbol{\psi}^{-1}$, so that

$$\mathbf{X} = \boldsymbol{\psi}(\mathbf{Z}; \mathbf{u}) \quad \text{and} \quad \mathbf{Z} = \boldsymbol{\phi}(\mathbf{X}; \mathbf{u}),$$

where \mathbf{u} is a vector representing all, if any, parameters. If, under the target density, $\boldsymbol{\phi}$ maps \mathbf{X} to a random object \mathbf{Z} that has a simple *base distribution*, such as a (multivariate) normal or uniform distribution, then the target distribution can be sampled indirectly by sampling from the base distribution and mapping the result using $\boldsymbol{\psi}$. The name of such a generative model is derived from the idea that a sequence of invertible differentiable transformations can map even a very complex probability distribution to a simple base distribution, thereby forming a ‘normalizing flow’. The construction of the ‘flow’ relies on the principle that any composition of invertible differentiable functions will also be invertible and differentiable.

In what follows, we assume that \mathcal{X} is a subset of \mathbb{R}^n and that \mathbf{X} and \mathbf{Z} are ‘‘continuous’’ random variables; more precisely, that they have probability densities $p_{\mathbf{X}}$ and $p_{\mathbf{Z}}$ with respect to the Lebesgue measure on \mathcal{X} . Since $\boldsymbol{\psi}$ is differentiable, the probability density function \mathbf{X} can be expressed in terms of the probability density function \mathbf{Z} via the relation

$$p_{\mathbf{X}}(\mathbf{x}) = p_{\mathbf{Z}}(\boldsymbol{\phi}(\mathbf{x}; \mathbf{u})) |\det \mathbf{D}\boldsymbol{\phi}(\mathbf{x}; \mathbf{u})|, \quad (8)$$

or, in terms of \mathbf{z} :

$$p_{\mathbf{X}}(\boldsymbol{\psi}(\mathbf{z}; \mathbf{u})) = p_{\mathbf{Z}}(\mathbf{z}) |\det \mathbf{D}\boldsymbol{\psi}(\mathbf{z}; \mathbf{u})|^{-1}. \quad (9)$$

Here, $\mathbf{D}\boldsymbol{\phi}(\mathbf{x}; \mathbf{u})$ is Jacobian matrix of $\boldsymbol{\phi}$ and the absolute value of its determinant is the *Jacobian* of $\boldsymbol{\phi}$, and similar for $\boldsymbol{\psi}$. If $\boldsymbol{\psi}$ is a composition of other invertible differentiable functions,

$$\boldsymbol{\psi}(\cdot; \mathbf{u}) = \circ_i \boldsymbol{\psi}_i(\cdot; \mathbf{u}_i),$$

then we can use the chain rule to combine the Jacobians as a product,

$$|\det \mathbf{D}\boldsymbol{\psi}(\cdot; \mathbf{u})| = \prod_i |\det \mathbf{D}\boldsymbol{\psi}_i(\cdot; \mathbf{u}_i)|.$$

Therefore, analogous to a feed-forward neural network, the full Jacobian can be computed during a single pass as $\boldsymbol{\psi}(\mathbf{z}; \mathbf{u})$ is computed. A similar result holds for the inverse: $\boldsymbol{\phi}$. Functional composition increases the complexity of the model while maintaining the accessibility of the Jacobian. In this way, Normalizing Flows provides the opportunity for highly expressible generative models, with computationally tractable density functions. As will be discussed in Section 2.4, training the model to learn a target density from a given function, rather than via training data, requires the ability to compute $p_{\mathbf{X}}$ as a differentiable function. Therefore, using a Normalizing Flows model improves and simplifies the framework introduced by Gibson and Kroese (2022) by eliminating the need to train a second neural network to approximate the density function via kernel density estimation.

2.3 Coupling Flows

Coupling Flows, first introduced by Dinh, Krueger, and Bengio (2015), is a family of Normalizing Flows which have been shown to be universal approximators of arbitrary invertible differentiable functions when composed appropriately (Teshima et al., 2020). A single Coupling Flows unit splits \mathbf{z} into two vectors, \mathbf{z}^A and \mathbf{z}^B , as well as \mathbf{x} into \mathbf{x}^A and \mathbf{x}^B . It then transforms one part via an invertible differentiable function, \boldsymbol{w} , called the *coupling function*, which contains parameters determined by the other partition. The relationship between the second part and the coupling function parameters, called the *conditioner*, Θ , does not need to be invertible, since this part is not transformed and is accessible when computing the inverse. This means functions with many learnable parameters, like neural networks, can be used to capture interesting and complex dependencies between variables. Figure 1 illustrates the structure of a Coupling Flow models and its inverse.

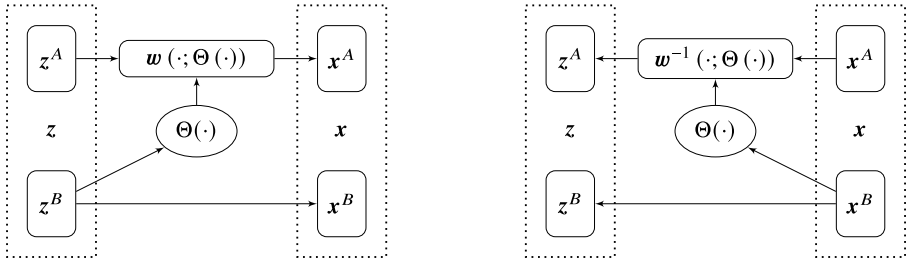


Fig. 1 Flow charts showing the structure of a Coupling Flow (left) and its inverse (right).

The Coupling Flow can be expressed mathematically as:

$$\begin{aligned} \mathbf{x} &= (\mathbf{x}^A, \mathbf{x}^B) = \boldsymbol{\psi}(\mathbf{z}^A, \mathbf{z}^B; \mathbf{u}) = (\mathbf{w}(\mathbf{z}^A; \Theta(\mathbf{z}^B; \mathbf{u})), \mathbf{z}^B), \\ \mathbf{z} &= (\mathbf{z}^A, \mathbf{z}^B) = \boldsymbol{\phi}(\mathbf{x}^A, \mathbf{x}^B; \mathbf{u}) = (\mathbf{w}^{-1}(\mathbf{z}^A; \Theta(\mathbf{x}^B; \mathbf{u})), \mathbf{x}^B). \end{aligned}$$

Note that the parameters \mathbf{u} of the flow are contained in the conditioner, Θ . Since the B -partition does not transform, the Jacobian of the full transformation is just the Jacobian of the coupling function; that is,

$$|\det D\boldsymbol{\psi}(\mathbf{z}, \mathbf{u})| = |\det D\mathbf{w}(\mathbf{z}^A; \Theta(\mathbf{z}^B; \mathbf{u}))|.$$

The coupling function is often chosen to be a simple affine transformation (Dinh et al., 2015; Dinh, Sohl-Dickstein, & Bengio, 2016; Kingma & Dhariwal, 2018). Instead, we base our coupling function on the rational function proposed by Ziegler and Rush (2019):

$$r(z; \theta_1, \theta_2, \theta_3, \theta_4, \theta_5) = \theta_1 z + \theta_2 + \frac{\theta_3}{1 + (\theta_4 z + \theta_5)^2}. \quad (10)$$

This function has five parameters when applied elementwise. However, separate parameters can be used for each dimension in the \mathbf{z}^A partition. For additional complexity, the coupling function can also be a composition of these rational functions, $w(\cdot; \Theta) = \circ_i r_i(\cdot; \theta_1^i, \theta_2^i, \theta_3^i, \theta_4^i, \theta_5^i)$. In this case the total number of parameters in the coupling function is five multiplied by the number of compositions multiplied by the number of dimension in \mathbf{z}^A . We choose the conditioner function to be a Multilayer Perceptron, with an input dimensionality to match \mathbf{z}^B and an output dimensionality to match the number of parameters in the coupling function.

In order for eq. (10) to define an invertible function, some restrictions have to be placed on the parameters. Choosing

$$\theta_1 > 0, \quad \theta_4 > 0, \quad \text{and} \quad |\theta_3| < \frac{8\sqrt{3}\theta_1}{9\theta_4}$$

ensures that the function r in eq. (10) is a strictly increasing invertible function. We enforce these constraints, following Ziegler and Rush, by processing the output of the conditioner Multilayer Perceptron,

$$\theta_1 = e^{\theta'_1}, \quad \theta_2 = \theta'_2, \quad \theta_4 = e^{\theta'_4}, \quad \theta_5 = \theta'_5, \quad \theta_3 = 0.95 \frac{8\sqrt{3}\theta_2}{9\theta_4} \tanh \theta'_3,$$

where primed letters denote unconstrained outputs of the conditioner function. The factor of 0.95 in the expression for θ_3 helps improve stability by precluding barely invertible functions that can exist near the bound $|\theta_3| = 8\sqrt{3}\theta_1/(9\theta_4)$. While we do not need to compute the inverse in our method, it can be calculated as the real solution to a cubic equation.

A single coupling flow unit does not transform the \mathcal{B} -partition. So, functional composition is necessary to obtain a general approximation. To achieve this, the dimensions of each coupling flow unit are split differently to ensure that all dimensions have to opportunity to ‘influence’ all other dimensions, thereby capturing any dependencies between all dimensions. In practice, we permute the dimensions between Coupling Flow units, while fixing the partitioned dimensions. The permutation of dimensions can be viewed as special case of a linear transformation, where the transformation matrix is a permutation of the rows of the identity matrix. The Jacobian matrix of a linear transformation is equal to the transformation matrix. So, the Jacobian is just 1 and is trivial to include in the probability density calculations of eq. (8) and eq. (9).

2.4 Objective and Training

Choosing an appropriate Normalizing Flows model architecture allows any density to be approximated arbitrarily well if the correct parameters can be identified. We use a form of stochastic gradient descent to iteratively tune the parameters of the model to minimize the Kullback–Leibler (KL) divergence between the model density, q say, and a known target function, $h : \mathcal{X} \rightarrow \mathbb{R}$, that is proportional to the target probability density function. That is, we minimize

$$\mathcal{D}(q, h) := \mathbb{E}_q \ln \frac{q(\mathbf{X})}{h(\mathbf{X})}. \quad (11)$$

Note that a model density q that minimizes eq. (11) also minimizes the KL divergence between q and the actual (normalized) target density. The advantage of using eq. (11) is that the normalization constant of the target density does not need to be known. This motivates the following minimization program:¹

$$\min_{\mathbf{u}} L(\mathbf{u}) := \min_{\mathbf{u}} \mathbb{E} [\ln p_{\mathbf{Z}}(\mathbf{Z}) - \ln |\det D\psi(\mathbf{Z}; \mathbf{u})| - \ln h(\psi(\mathbf{Z}; \mathbf{u}))], \quad (12)$$

¹Actually, the first term in this objective corresponds to the negative differential entropy of the base density and does not depend on any parameters, so could be omitted. However, we keep it to make interpreting the loss a little simpler.

which corresponds to choosing parameters to minimize the KL divergence between the Normalizing Flows generative density and a target density whose probability density function is proportional to h . The expression can be derived by substituting q from eq. (9) into eq. (11) and taking the expectation with respect to the base density. The minimum KL divergence is zero, when the two densities are identical, so the minimum value of this objective is the negative natural logarithm of the normalization constant of h . For example, if $h(\mathbf{x}) = p(\mathbf{x})H(\mathbf{x})$, then the minimum value of the objective function L is $-\ln \ell$. The objective function values can be estimated via

$$\widehat{L}(\mathbf{u}) := \frac{1}{n} \sum_{k=1}^n [\ln p_{\mathbf{Z}}(\mathbf{Z}_k) - \ln |\det \mathbf{D}\boldsymbol{\psi}(\mathbf{Z}_k; \mathbf{u})| - \ln h(\boldsymbol{\psi}(\mathbf{Z}_k; \mathbf{u}))], \quad (13)$$

where $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ is an iid sample of the base density.

The objective function L marks a clear distinction between how we train a Normalizing Flows generative model using a target function, and how they are typically trained using a data set. In other works, such as by Dinh et al. (2015), Normalizing Flows are trained to learn the density of a data set by maximizing the log-likelihood of the sampled data. While this can be considered equivalent to minimizing the KL divergence, the training process involves sampling the data set, which corresponds to sampling \mathbf{X} , and then moving in the ‘normalizing’ direction with $\boldsymbol{\phi}$ to solve the maximization program:

$$\max_{\mathbf{u}} \mathbb{E}_q \ln q(\mathbf{X}) = \max_{\mathbf{u}} \mathbb{E}_q [\ln p_{\mathbf{Z}}(\boldsymbol{\phi}(\mathbf{X}; \mathbf{u})) + \ln |\det \mathbf{D}\boldsymbol{\phi}(\mathbf{X}; \mathbf{u})|].$$

After the model is trained, new data can be generated by sampling \mathbf{Z} and moving in the ‘generating’ direction using $\mathbf{X} = \boldsymbol{\psi}(\mathbf{Z}; \mathbf{u})$. Therefore, in the data-focused case, it is essential to be able to compute both transformations, $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$. However, in the context of rare-event simulation, we do not have access to training data, but a target function that is proportional to a target probability density, and all training occurs in the ‘generating’ direction. Therefore, computing $\boldsymbol{\phi}$ is not required (although $\boldsymbol{\psi}$ must still be invertible in principle). Algorithm 1 outlines our training procedure.

Algorithm 1 Training a flow-based generative model using a target function.

- 1: Randomly initialize parameters, \mathbf{u} .
 - 2: **for** many iterations **do**
 - 3: Independently sample $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ from the base distribution.
 - 4: Compute $\ln |\det \mathbf{D}\boldsymbol{\psi}(\mathbf{Z}_k; \mathbf{u})|$ and $\mathbf{X}_k = \boldsymbol{\psi}(\mathbf{Z}_k; \mathbf{u})$ during a single pass.
 - 5: Compute the target function values $h(\mathbf{X}_k)$.
 - 6: Estimate the objective function via eq. (13).
 - 7: Adjust parameters \mathbf{u} via standard gradient descent methods.
 - 8: **end for**
-

2.5 Conditional Target Densities

As previously shown in eq. (4) and eq. (7), when estimating rare-event probabilities and other quantity conditioned on a rare event, the target density includes an indicator function, $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$. However, the term, $-\ln h(\boldsymbol{\psi}(\mathbf{Z}; \mathbf{u}))$ in the objective function L in eq. (12) requires each value $h(\mathbf{x})$ of the target function value to be strictly positive. To circumvent this problem, we reuse the solution presented by Gibson and Kroese (2022), to approximate $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$ by defining a strictly positive *penalty factor*,

$$\rho(\mathbf{x}) := \exp[-\alpha(\gamma - S(\mathbf{x})) \mathbb{1}_{A^c}(\mathbf{x})], \quad (14)$$

which is a positive approximation of $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$ when $\alpha > 0$ and $\mathcal{A}^c = \{\mathbf{x} \in \mathcal{X} \mid S(\mathbf{x}) < \gamma\}$ is the complement set of \mathcal{A} . The approximation is more accurate the larger the value of α and is exact in the limit $\alpha \rightarrow \infty$. Figure 2 illustrates the convergence of the approximation.

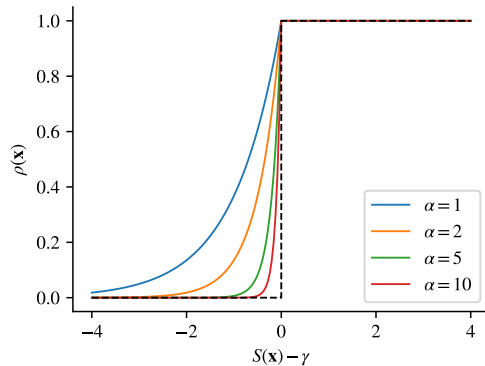


Fig. 2 A comparison between the penalty factor $\rho(\mathbf{x})$ (colored lines) and the step function $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$ (black dashed line). As α grows the approximation becomes more accurate.

Therefore, if the target density includes the indicator function as a factor, we replace it with the penalty factor with appropriate performance function, S . For example, if the aim is to estimate $\ell^{\mathcal{A}}$ from eq. (5), then the target function would be of the form

$$h(\mathbf{x}) = p(\mathbf{x})H(\mathbf{x})\rho(\mathbf{x}).$$

Note that in this case

$$\mathbb{E}[-\ln h(\mathbf{X})] = \mathbb{E}[-\ln p(\mathbf{X})H(\mathbf{X})] - \alpha \mathbb{E}[(\gamma - S(\mathbf{X})) \mathbb{1}_{A^c}(\mathbf{X})],$$

so learning a conditional density by including the penalty factor in the target function is equivalent to learning the unconditional target $p(\mathbf{x})H(\mathbf{x})$ and

including an additional penalty term in the objective that penalizes samples outside the rare-event region. In this way, α can be interpreted as the weight of the penalty. For our experiments we choose $\alpha = 100$.

3 Results

In this section we look at how well Normalizing Flows models can learn various target functions using several practical examples.

3.1 Truncated Normal Density

Firstly we consider a very simple one-dimensional example of a truncated normal density. Let X have a standard normal distribution: $X \sim \mathcal{N}(0, 1)$ and consider the rare-event region $\mathcal{A} = [\gamma, \infty)$. An obvious choice for the performance function is $\mathcal{S}(x) = x$, so if the goal is to estimate the rare-event probability $c = \mathbb{P}(X \geq \gamma)$, then the target function is

$$h(x) := p_X(x) \exp \left[-\alpha (\gamma - x) \mathbb{1}_{(-\infty, \gamma)}(x) \right].$$

A Coupling Flow is unsuitable for a one-dimensional problem, so we use a composition of three rational functions from eq. (10), containing a total of 15 parameters. Choosing a batch size of $n = 1,000$, a learning rate of 0.001, a weight decay parameter of 0.0001 and $\alpha = 100$, the model was trained via Algorithm 1 and the Adam gradient descent optimizer (Kingma & Ba, 2017) for 30,000 iterations. Figure 3 illustrates how the model converges towards the truncated normal target density, with threshold parameter $\gamma = 3$.

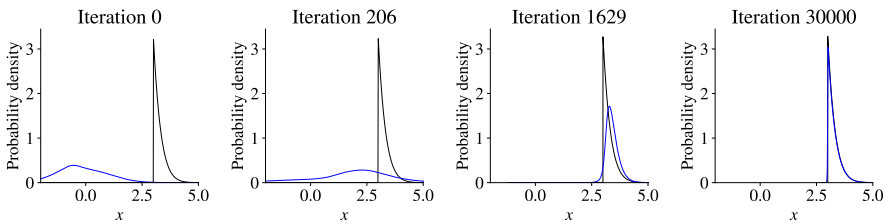


Fig. 3 The normalizing flow probability density function at different stages of training. The black density is the target probability density of a normal density truncated to the interval $[3, \infty)$. The blue density is the probability density of X at the corresponding training iteration.

Using a sample of 1,000 points, eq. (3) gives us an estimate 0.00134576 of the rare-event probability $\mathbb{P}(X \geq 3)$, which is about 0.3% error from the true value of about 0.0013499. From the same sample, the sample standard deviation is about 0.000261. We can estimate the minimum sample size needed

for a 1% standard error by

$$n = \left(\frac{\sigma}{\hat{c}} \frac{1}{0.01} \right)^2 \approx 376,$$

where σ is the sample standard deviation of the summand in eq. (3). Using a crude Monte Carlo estimator, the relative standard error should be $\sqrt{c(1-c)} \approx 0.0367$, requiring a much larger sample size of about $n = 7.4 \times 10^6$ for a 1% standard error.

The learned density is not exactly identical to the target conditional density, but is very close. To quantify how good the approximation is, we can estimate the KL divergence between the target density q^* and the learned density using

$$\mathcal{D}(q^*, q) = \mathbb{E}_{q^*} \left[\ln \frac{q^*(\mathbf{X})}{q(\mathbf{X})} \right] = \mathbb{E}_q \left[\frac{q^*(\mathbf{X})}{q(\mathbf{X})} \ln \frac{q^*(\mathbf{X})}{q(\mathbf{X})} \right].$$

In this example the KL divergence is estimated to be 0.03465 with a standard error of 0.56% using sample size 10,000.

3.2 Two-Dimensional Exponential Density

Secondly, we try another example from [Gibson and Kroese \(2022\)](#) of a two-dimensional exponential density. In particular, in this example the probability density function of $\mathbf{X} = (X_1, X_2)$ is:

$$p(x_1, x_2) = e^{-(x_1+x_2)},$$

and the rare-event region is $\mathcal{A} = \{(x_1, x_2) \in \mathbb{R}_+^2 : x_1+x_2 \geq \gamma\}$, with performance function, $S(x_1, x_2) = x_1+x_2$. To estimate $c = \mathbb{P}(X_1+X_2 \geq \gamma)$, the target function is

$$h(x_1, x_2) = \exp \left[-(x_1 + x_2) - \alpha (\gamma - x_1 - x_2) \mathbb{1}_{\{x_1+x_2 < \gamma\}} \right].$$

This time we choose a flow model composed of six Coupling Flows, interlaced with dimension permutations, followed by an elementwise exponential function. The exponentiation ensures that the generated points are positive. The coupling function in each Coupling Flow is a composition of two rational functions of the form in eq. (10). The conditioner in each Coupling Flow is a linear transformation plus a constant vector. Therefore, the total number of learnable parameters is 60. Figure 4 shows the structure of this model via a flow chart.

The model was trained for 100,000 iterations with a learning rate of 0.0001, a weight decay of constant of 0.0001, and a batch size of 10,000. When $\gamma = 10$, $c = (1 + \gamma)e^{-\gamma} \approx 0.000499399$. After training, with a sample size of 1,000 the estimated constant is 0.00049920, with standard deviation of about 6.34×10^{-5} giving a relative standard error of 0.40%. Achieving a comparable standard error using a crude Mote Carlo estimator would require a sample size of more than 10^8 . Figure 5 shows how well the model learned the target density by

including a scatter plot of a sample of points as well as comparing the learned probability density with the target probability density. The KL-divergence is estimated to be about 0.011 with 9.7% standard error using sample size of 10,000.

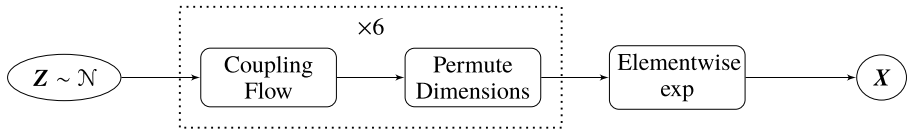


Fig. 4 Flow chart showing the normalizing flow structure used in the exponential density example. The base density is a 2D multivariate normal density which forms the input to a sequence of six consecutive Coupling Flow and dimension permutation pairs, followed by an exponential function applied to both dimensions. The learnable parameters are found in the conditioner functions of all six Coupling Flows.

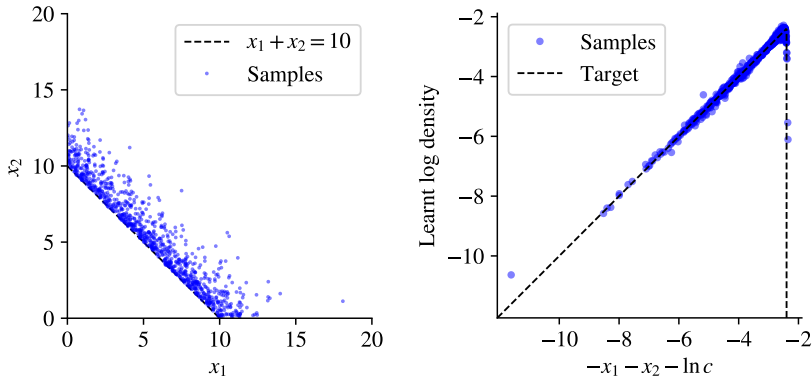


Fig. 5 1,000 points sampled from the trained normalizing flow. The left scatter plot shows the how the sampled two-dimensional points almost exclusively satisfy the rare-event condition, $x_1 + x_2 \geq 10$. The right plot compares the learned log-density with the renormalized exponential log-density. The dashed black line shows the target density, which is equal to the renormalized log density, until about -2.398 , which is the maximum log-density when $x_1 + x_2 = 10$.

3.3 Bridge Network

This third example comes from Chapter 9 of Kroese et al. (2011). In this example, we increase the number of dimensions to five, each corresponding to a random edge length in a bridge network, shown in Figure 6. The graph contains five edges and four nodes, where the length of each edge is uniformly distributed and the goal is to identify the expected shortest path from node A to node D . Specifically, we wish to estimate the expected length of the shortest path from A to D ; that is, the expectation ℓ in eq. (1), with $X_i \sim \mathcal{U}(0, 1)$ and

$$H(X) = \min\{X_1 + X_4, X_1 + 3X_3 + 2X_5, 2X_2 + 3X_3 + X_4, 2X_2 + 2X_5\}.$$

In this case $p(\mathbf{x}) = 1$, so the optimal Importance Sampling density that minimizes the variance in estimating ℓ is $q^*(\mathbf{x}) = H(\mathbf{x})/\ell$. Therefore, the target function is $h(\mathbf{x}) = H(\mathbf{x})$.

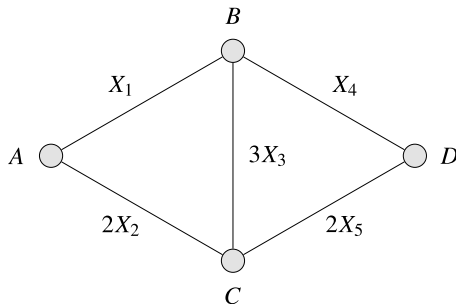


Fig. 6 A bridge network containing four nodes and five edges, where the edge lengths X_1, \dots, X_5 are random variables. So the minimum path length from node A to node D is also a random variable.

The space of allowed values in this example is $\mathcal{X} = [0, 1]^5$, so it is reasonable to choose a model architecture in which the domains and co-domains of each composed flow unit is $[0, 1]^5$. To achieve this we choose a uniform base distribution, $Z_i \sim \mathcal{U}(0, 1)$. Additionally, the model is composed of five Coupling Flows and dimension permutation pairs. In each permutation the dimensions are cycled via $[3, 4, 5, 1, 2]$. The dimensions are partitioned so the coupling function transforms the first three dimensions. The coupling function of each Coupling Flow is a rational function based on eq. (10) with separate parameters for each dimension. The conditioner, like the previous example, is a linear transformation. Figure 7 uses a flow chart to illustrate the structure of this model. The parameters in the coupling functions of each Coupling Flow unit contain additional constraints to ensure that $r(0) = 0$ and $r(1) = 1$, thereby preserving the domain of $[0, 1]^5$ at each flow component.

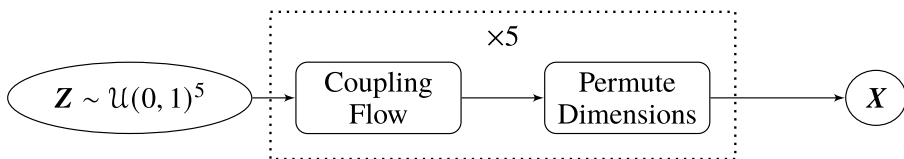


Fig. 7 Flow chart showing the normalizing flow structure used in the bridge network example. The base distribution is a 5D uniform distribution which forms the input to a sequence of five consecutive Coupling Flows and dimension permutation pairs.

The model was trained with a learning rate of 0.0001, weight decay parameter of 0.0001, batch size 10,000 for 300,000 iterations. The scatter plot in Figure 8 compares the learned probability density of $\psi(\mathbf{Z}; \mathbf{u})$ with the optimal sampling density given by eq. (2). The learned density forms a fairly good

approximation of the optimal sampling density with a KL divergence of about 0.0017 with 29% standard error (using sample of 10,000). The histograms in Figure 8 compare the densities of the summands of the Importance Sampling estimator from eq. (1) and the crude Monte Carlo estimator, being the sample mean of $H(\mathbf{X})$. In this example both methods provide accurate estimates of the expected minimum path length. Typical values of the crude Monte Carlo and IS estimators are 0.9272 with a standard relative error of 0.43% and 0.9301 with a standard relative error of 0.053%. While both values agree with the theoretical value of $l = \frac{1339}{1440} = 0.92986\bar{1}$, the IS estimator reduces the variance by a factor of about 65.

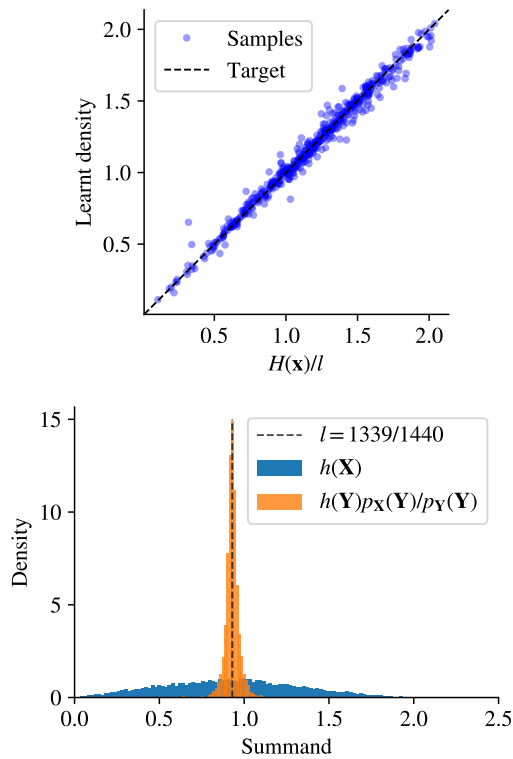


Fig. 8 A scatter plot (left) comparing the learned probability density and the target density, and a histogram (right) comparing the summand densities of the crude Monte Carlo estimator (blue) and the Importance Sampling estimator (orange). Clearly the model has learned to approximate the target density and can estimate the expected minimum path length with a much lower variance using Importance Sampling than the crude Monte Carlo estimator.

3.4 Conditional Bridge Network

In this example we continue with the bridge network, but now consider the rare event that the shortest path includes three edges, passing through the

middle, rather than just two. That is, the shortest path is either $ABCD$ or $ACBD$. The rare-event region can be defined by choosing

$$S(\mathbf{X}) = \min(X_1 + 3X_3 + 2X_5, 2X_2 + 3X_3 + X_4) - \min(X_1 + X_4, 2X_2 + 2X_5),$$

and $\gamma = 0$. That is, the penalty depends on the difference between the shortest ‘inner’ path ($ABCD$ or $ACBD$) and the shortest ‘outer’ path (ABD or ACD).

The exact same model was reused, but trained across 500,000 iterations with the penalty factor included in the target function. In this way, the goal is to estimate the expected shortest path length of the distribution conditioned on the shortest path being $ABCD$ or $ACBD$. The rare event probability can be estimated using eq. (3) with a sample size of 10,000 to be about 0.0346 with a standard relative error of 1.3%. This corresponds to a variance reduction in a factor of about 17 compared to the crude Monte Carlo estimate. The conditional expected shortest path length can now be estimated via eq. (6) as 0.913 with a standard relative error of 1.7%.

By learning an approximation of the conditional distribution, we can explore the likely conditions required to generate the rare event. Figure 9 compares scatter plots between the learned unconditional and the conditional distributions. The unconditional distribution is not too different from the original uniform distribution, with the exception of X_1 and X_4 . This is the expected result since the path, ABD , with path length $X_1 + X_4$, is the most likely to be the shortest. However, the learned conditional distribution is significantly different to the original uniform distribution. From the scatter plots it is clear that X_3 is almost always short, and X_2 and X_5 have a strong negative correlation. Therefore, we can conclude that the most likely conditions for the shortest path to be one that passes through the middle edge, is that the middle edge is short, and that exactly one of bottom two edges is short. This is congruent with expectations since a long middle edge would likely make the total path length longer than both of the outer paths, if both bottom edges were short then the middle edge would be bypassed, and if both bottom edges were long, then the shortest path would likely be across the top.

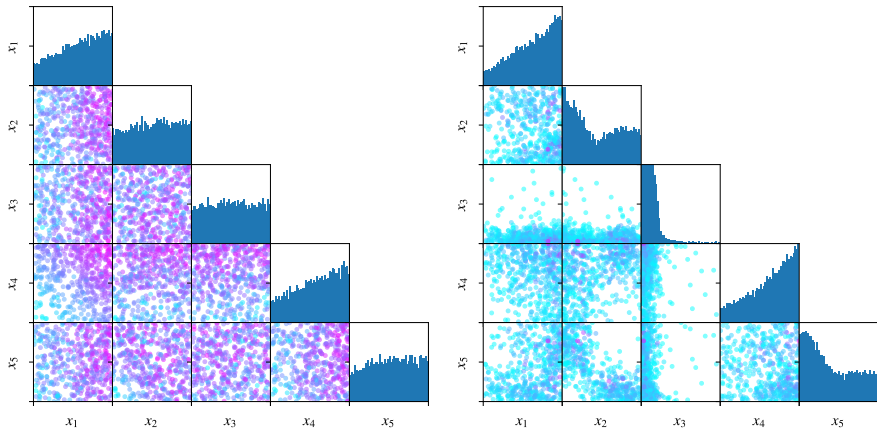


Fig. 9 Scatter plots comparing each of the five dimensions against each other, as well as histograms showing the distribution of each dimension. The darker purple colour corresponds to a higher probability density. The left figure shows the learned unconditional distribution while the right shows the learned conditional distribution.

3.5 Asian Call Option

In this section we look at stock price models and expected pricing of call options. This example, from Chapter 15 of Kroese et al. (2011), specifically looks at estimating the expected Asian call option of a stock whose undiscounted price, S_t , at time t can be modelled by the Stochastic Differential Equation (SDE)

$$dS_t = rS_t dt + \sigma S_t dX_t,$$

where the drift, r , is the risk-free rate of return, $\{X_t\}$ is a Wiener process and σ is the volatility. This SDE is solved by

$$S_t = S_0 e^{(r - \sigma^2/2)t + \sigma X_t}.$$

The average stock price across the interval $t \in [0, T]$ can be approximated by

$$\bar{S}_T = \frac{1}{n+1} \sum_{i=0}^n S_{t_i}, \quad t_i = \frac{iT}{n}, \quad i = 0, \dots, n,$$

where the time has been discretized into n equal intervals of $\Delta t = T/n$,

$$S_{t_i} = S_0 \exp\left((r - \sigma^2/2)t_i + \sigma X_{t_i}\right),$$

and $X_{t_i} \sim \mathcal{N}(0, \Delta t)$. The discounted payoff at time $t = 0$ of the Asian option is

$$H(\mathbf{X}) = e^{-rT} (\bar{S}_T - K)^+,$$

where $\mathbf{X} = (X_{t_1}, \dots, X_{t_n})$ and the $+$ superscript denotes the ramp function. Since $\{X_t\}$ is a Wiener process, $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \Delta t I)$, where I is the identity matrix of appropriate dimensionality. If the goal is to estimate the expected discounted payoff then, according to eq. (2), the optimal sampling distribution should include a factor of $H(\mathbf{x})$. However, the objective eq. (12) requires the target function to be non-zero at all sampled values. Therefore, we introduce the following approximation

$$x^+ \approx v(x) := \begin{cases} x, & \text{if } x \geq \delta \\ \delta e^{\frac{x}{\delta}-1}, & \text{if } x < \delta \end{cases}$$

which is a strictly positive continuously differentiable function that is exact in the limit $\delta \rightarrow 0$. In our experiments we choose $\delta = 0.5$. Therefore, the target function is

$$h(\mathbf{x}) = v\left(e^{-rT}(\bar{S}_T - K)\right) \frac{1}{(2\pi\Delta t)^{n/2}} e^{-\frac{1}{2\Delta t}|\mathbf{x}|^2}.$$

We chose the base distribution to be $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \Delta t I)$ and the normalizing flow to be comprised of six coupling flow units with equal partition sizes. The dimensions were permuted after each coupling flow unit. Figure 10 illustrates the normalizing flows structure.

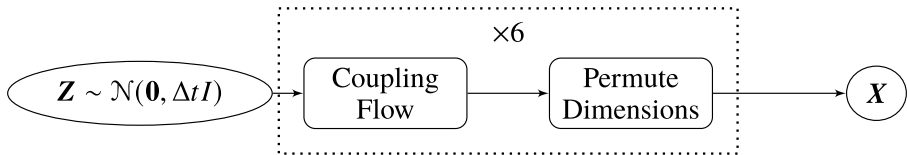


Fig. 10 Flow chart showing the normalizing flow structure used in the Asian call option example. The base density is a 88D multivariate normal density which forms the input to a sequence of six consecutive Coupling Flow and dimension permutation pairs.

The parameters used in this example were $r = 0.07$, $\sigma = 0.2$, $K = 35$, $S_0 = 40$, $T = 4/12$ and $n = 88$. The model was trained for 30000 steps with a batch size of 1000, learning rate of 0.0001 and weight decay of 0.0001. Figure 11 compares the histograms of the crude Monte Carlo and importance sampling estimator summands. Using a sample of 10000, the crude estimator is about 5.3432 with a standard relative error of 0.49%, and the importance sampling estimator is about 5.3580 with a standard relative error of 0.23%. This corresponds to a variance reduction of about a factor of 4.4.

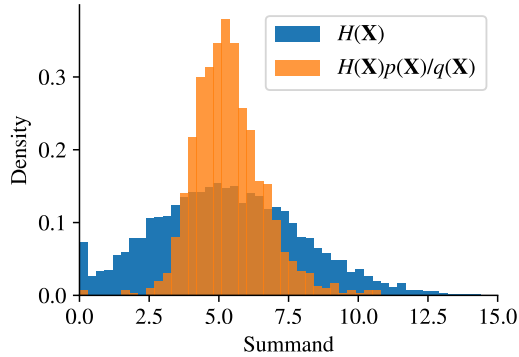


Fig. 11 Histograms of the crude Monte Carlo and Importance Sampling estimator summands. Importance Sampling using the learned distribution allows the expected discounted payoff of the Asian option to be estimated with less variance than the crude Monte Carlo estimate.

Let us explore the rare event in which the discounted payoff exceeds 14, $\mathcal{A} = \{\mathbf{x} \in \mathbb{R}^n : e^{-rT}(\bar{S}_T - K)^+ \geq \gamma\}$ with performance function $S(\mathbf{x}) = e^{-rT}(\bar{S}_T - K) \geq \gamma$ and $\gamma = 14$. In this case the goal is to explore how this rare event occurs and to estimate the probability of the rare event. Therefore, the target function is just the probability density function of the Wiener process multiplied by the penalty factor,

$$h(\mathbf{x}) = \frac{1}{(2\pi\Delta t)^{n/2}} e^{-\frac{1}{2\Delta t}|\mathbf{x}|^2} \rho(\mathbf{x}).$$

This time the model was trained for 100000 iterations. With a sample of 10000 the crude and IS estimators are about 0.0022 with 21% relative error, and 0.0015797 with 0.46% relative error respectively. This corresponds to a variance reduction of a factor of about 2100. Using the estimated normalization constant of 0.0015797 the KL divergence between the learnt density and the target density is about 0.15157 with a standard relative error of 0.40% indicating that the learnt distribution is very close to the target conditional distribution. Figure 12 illustrates the strong linear relationship between the learnt log density and the target unconditional density, and compares the distributions of the simulated discounted Asian options.

Since the learnt distribution is a decent approximation of the conditional distribution, we can use the simulated trajectories to gain insights into how high performing options can occur. Figure 13 shows typical stock prices across time when simulating using the original probability measure and the learnt probability measure. In the first instance the average value is largely static, increasing very slowly, while the variance grows across time. However, conditioning the simulation on achieving a discounted payoff of at least 14 changes the trend unexpectedly. Rather than a steady growth in stock price across the whole time interval, most of the stock price growth occurs by $t = 0.25$.

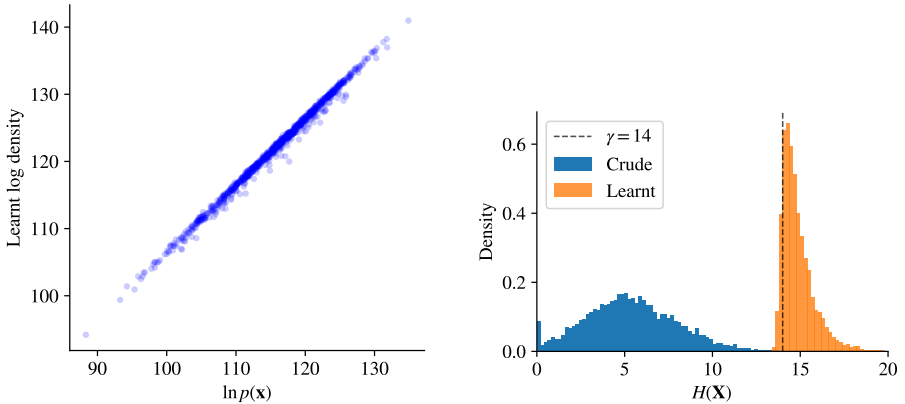


Fig. 12 A scatter plot comparing the learnt log probability density and the unconditional log probability density (left) and Histograms (right) of simulated discounted payoffs using crude Monte Carlo (blue) and the learned conditional distribution (orange).

Additionally, the variance is fairly constant over time during this growth time, mostly growing near the end. This result is also visible in the covariance matrix of the stock price, illustrated in figure 14. Under the original probability distribution the stock price near maturity depends little on the early prices, and the variance grows roughly linearly with time. However, under the learnt probability distribution, the variance in stock price is largely constant until a short time before maturity. Additionally, there is a small correlation between early stock prices and late stock prices. This could indicate that

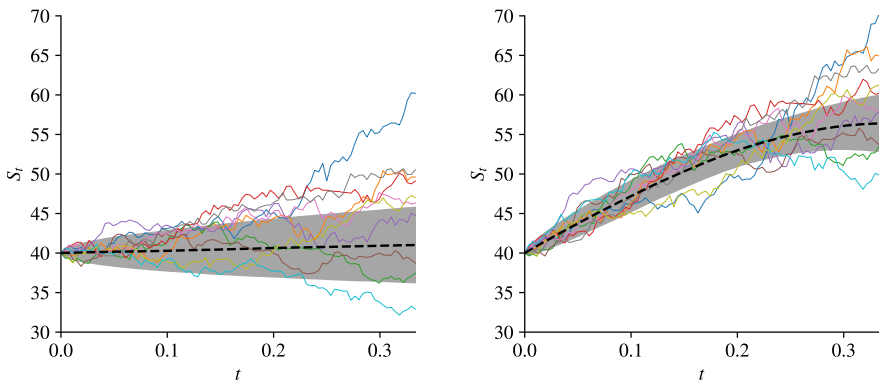


Fig. 13 Undiscounted stock prices generated by crude Monte Carlo (left) and the model conditioned on discounted payoffs of at least 14 (right). Coloured lines represent 10 trajectories, dashed black lines represent the mean and the shaded grey region is a single standard deviation from the mean.

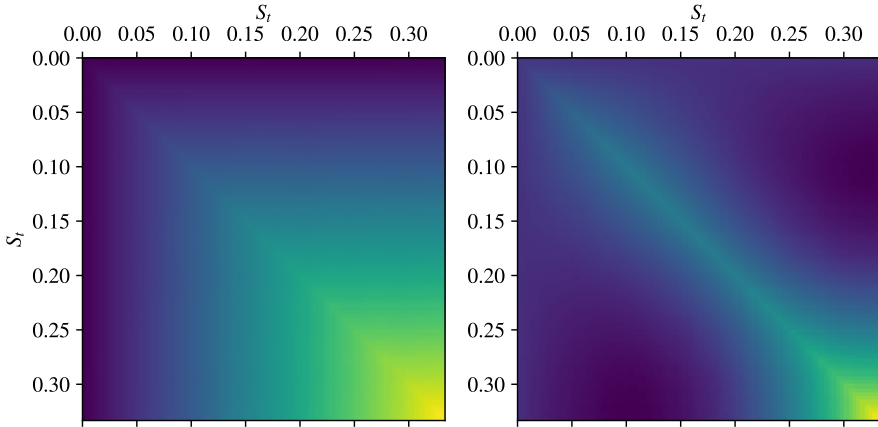


Fig. 14 Covariance matrix of undiscounted stock prices simulated using crude Monte Carlo (left) and generated by the model conditioned on discounted payoffs of at least 14 (right). In the crude simulation the variance grows with with time. There is a small amount of covariance between early stock prices and later stock prices in the conditionally generated prices. It appears that while most stock prices achieve the discounted payoff of at least 14 through consistent steady increases, there are a few trajectories which perform well near the beginning or near the end, but not both.

3.6 Double Slit Experiment

In this example we test our method on a high-dimensional problem. A dimensionless particle traverses an unbounded 2-dimensional environment for a maximum duration of T , which is discretized into $\Delta t = \frac{d}{2}$ time steps, where $d \geq 2$. The environment consists of a barrier with two slits, located at $(x_{\text{slit}}, y_{\text{slit}})$ and $(x_{\text{slit}}, -y_{\text{slit}})$. The width of both slits is w_{slit} . Additionally, the environment contains a vertical screen located at x_{screen} . The aim of the particle is to reach the screen while simultaneously passing through one of the two slits. At time step 0, the particle starts at the $2D$ -position $\mathbf{V}_0 = (0, 0)$. Given a random vector $\mathbf{Q} \sim \mathcal{N}(\mathbf{0}, \frac{2T}{d}I_d)$, where I_d is the d -dimensional identity matrix, the position of the particle at time step $1 \leq k \leq d/2$ evolves according to

$$\mathbf{V}_k = \mathbf{V}_{k-1} + \mathbf{Q}_k, \quad (15)$$

where $\mathbf{Q}_k = (Q_{2k-1}, Q_{2k})^\top$ is the k -th component vector of \mathbf{Q} . The path \mathbf{V} of the particle induced by \mathbf{Q} is then the sequence of points $\mathbf{V} = (\mathbf{V}_k)_{k=1}^{d/2}$, with \mathbf{V}_k defined according to eq. (15).

Let g_d be the density of $\mathcal{N}(\mathbf{0}, \frac{2T}{d}I_d)$. Our goal is to learn a sampling density that is proportional to g_d truncated to the rare-event region $\mathcal{A} = \{\mathbf{q} \in \mathbb{R}^d : S(\mathbf{q}) \geq 0\}$, where \mathbf{q} is an outcome of \mathbf{Q} . The performance function $S(\mathbf{q})$ is defined as follows. Suppose is M is the time step where \mathbf{X} is closest to the screen; that is,

$$M := \operatorname{argmin}_{1 \leq k \leq d/2} (x_{\text{screen}} - X_k),$$

where X_k is the x -coordinate of the particle at time step k . Additionally, let N be the first step where $X_{N-1} \leq x_{\text{slit}} \leq X_N$, i.e., where the particle crosses x_{slit} . The performance function of an outcome \mathbf{q} of \mathcal{Q} , with path $\mathbf{v} = P(\mathbf{q})$, and outcomes n and m of N and M is then defined as

$$S(\mathbf{q}) = -\left(\min\{f_1(y_n), f_2(y_n)\} + f_3(x_m) \right), \quad (16)$$

where

$$\begin{aligned} f_1(y_n) &= \left(|y_n - y_{\text{slit}}| - \frac{w_{\text{slit}}}{2} \right) \mathbb{1}_{\{|y_n - y_{\text{slit}}| > \frac{w_{\text{slit}}}{2}\}}, \\ f_2(y_n) &= \left(|y_n + y_{\text{slit}}| - \frac{w_{\text{slit}}}{2} \right) \mathbb{1}_{\{|y_n + y_{\text{slit}}| > \frac{w_{\text{slit}}}{2}\}}, \\ f_3(x_m) &= (x_{\text{screen}} - x_m) \mathbb{1}_{\{x_{\text{screen}} - x_m > 0\}}. \end{aligned} \quad (17)$$

The functions f_1 and f_2 in eq. (17) measure the distances of the y -coordinate (i.e., y_n) of point \mathbf{v}_n to the slits and the min operator in eq. (16) ensures that only the minimum distance to the slits contributes to the performance function. Note that the definitions of f_1 and f_2 imply that a path successfully crosses a slit if $y_{\text{slit}} - \frac{w_{\text{slit}}}{2} \leq |y_n| \leq y_{\text{slit}} + \frac{w_{\text{slit}}}{2}$, even if the path touches the barrier from \mathbf{v}_{n-1} to \mathbf{v}_n . In case the particle does not cross x_{slit} , we set both f_1 and f_2 to 0. The function f_3 is an additional penalty term which penalizes the x -coordinate of the point \mathbf{v}_m (i.e., x_m) to be on the left-hand side of the screen.

With the performance function in eq. (16), the target function is defined as

$$h(\mathbf{q}) = g_d(\mathbf{q}) \exp \left[-\alpha(\gamma - S(\mathbf{q})) \mathbb{1}_{\{S(\mathbf{q}) < 0\}} \right], \quad (18)$$

with $\gamma = 0$. For the base distribution, we chose a Gaussian Mixture Model (GMM) consisting of two d -dimensional component distributions $\mathcal{N}(\boldsymbol{\mu}_1, \frac{2T}{d} I_d)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \frac{2T}{d} I_d)$ with equal weights, where $\boldsymbol{\mu}_1 = (1, -1, \dots, 1, -1)^\top$ and $\boldsymbol{\mu}_2 = (1, \dots, 1)^\top$. This bimodal mixture model helps in learning the correct bimodal target density. The normalizing flow is composed of twenty coupling flow units with equal partition sizes. The dimensions are permuted after each coupling flow unit. Figure 15 illustrates the normalizing flow.

In our experiments we set $T = 1$, $x_{\text{slit}} = 5$, $y_{\text{slit}} = 1.5$, $w_{\text{slit}} = 1$ and $x_{\text{screen}} = 10$. We then trained three normalizing flows using $d = 100$, $d = 500$ and $d = 1,000$, where each flow was trained for 100,000 iterations with a batch size of 200, learning rate of 10^{-6} and weight decay of 10^{-8} .

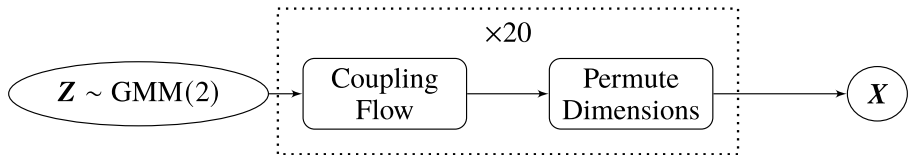


Fig. 15 Flow chart showing the normalizing flow structure used in the Double Slit example. The base density is a Gaussian Mixture Model (GMM) consisting of two component distributions $N(\mu_1, \frac{2T}{d} I_d)$ and $N(\mu_2, \frac{2T}{d} I_d)$ with equal weights, where $\mu_1 = (1, -1, \dots, 1, -1)^\top$ and $\mu_2 = (1, \dots, 1)^\top$. This forms the input to a sequence of twenty consecutive Coupling Flow and dimension permutation pairs.

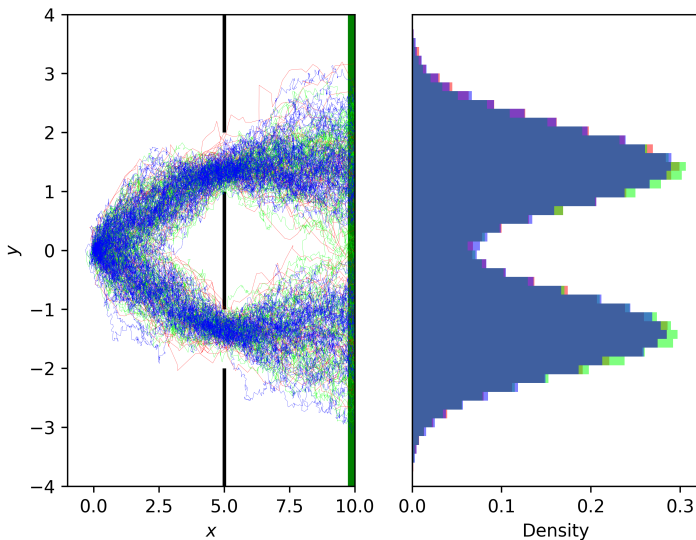


Fig. 16 The double slit environment (left figure) with 100 paths simulated with the learned normalizing flows for $d = 100$ (red), $d = 500$ (green) and $d = 1,000$ (blue). The black regions represent the barrier, while the dark green region represents the screen. The right figure shows the histograms of 100,000 simulated paths at the screen for $d = 100$ (red), $d = 500$ (green) and $d = 1,000$ (blue).

Figure 16 (left) shows the double slit environment with paths simulated using the learned normalizing flows for $d = 100$ (red paths), $d = 500$ (green paths) and $d = 1,000$ (blue paths). We additionally used 100,000 samples for each d to test the capability of the models in generating paths that successfully pass through one of the slits and hit the screen at x_{screen} . In our tests, the success rate was 91.6%, 91.2% and 90.3% for $d = 100$, $d = 500$ and $d = 1,000$ respectively, which shows that the learnt distributions are close to the conditional target distribution, even when the base and target distributions are 1,000-dimensional. We additionally investigate how the marginal distributions of the paths at x_{screen} look like. Figure 16 (right) shows the histogram of the

y-positions of 100,000 paths at x_{screen} for $d = 100$ (red), $d = 500$ (green) and $d = 1,000$ (blue) respectively. We can see that for all values of d , the distributions at the screen closely match each other. This is both desired and expected, due to the scaling factor $\frac{2T}{d}$ for the variance of the conditional target densities.

4 Conclusion

Monte Carlo methods are important tools for decision making under uncertainty. However, rare events pose significant challenges for standard Monte Carlo methods, due to the low probability of sampling such events. In this paper, we propose a framework for rare-event simulation. Our framework uses a normalizing flow based generative model to learn optimal importance sampling distributions, including conditional distributions. The model is trained using standard gradient descent techniques to minimize the KL divergence between the model density and a function that is proportional to the target density. Empirical evaluations in several domains demonstrate that our framework is able to closely approximate complicated distributions, even in high-dimensional (up to 1,000-dimensional) rare-event settings. Simultaneously, our framework demonstrates substantial improvements in sample efficiency compared to standard Monte Carlo methods. Given these properties, we believe that combining our framework with Monte Carlo based methods for decision making under uncertainty is a fruitful avenue for future research.

Supplementary information. The source-code of our framework and the examples is publicly available at <https://github.com/hoergems/rare-event-simulation-normalizing-flows>.

Acknowledgements. We thank Wei Jiang for helpful discussions. This work is supported by the Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS) CE140100049 grant and the Australian Research Council (ARC) Discovery Project 200101049.

References

- Ardizzone, L., Lüth, C., Kruse, J., Rother, C., Köthe, U. (2019). *Guided image generation with conditional invertible neural networks*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.1907.02392>
- Arief, M., Bai, Y., Ding, W., He, S., Huang, Z., Lam, H., Zhao, D. (2021). *Certifiable deep importance sampling for rare-event simulation of black-box systems*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.2111.02204>
- Botev, Z.I., Kroese, D.P., Taimre, T. (2007). Generalized cross-entropy methods with applications to rare-event simulation and optimization. *Simulation*, 83(11), 785–806, <https://doi.org/10.1177/0037549707087067>

- Bucklew, J. (2004). *Introduction to rare event simulation*. Springer.
- de Boer, P.-T., Kroese, D.P., Mannor, S., Rubinstein, R.Y. (2005). A tutorial on the cross-entropy method. *Annals of Operations Research*, 134(1), 19–67, <https://doi.org/10.1007/s10479-005-5724-z>
- Dinh, L., Krueger, D., Bengio, Y. (2015). *Nice: Non-linear independent components estimation*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.1410.8516>
- Dinh, L., Sohl-Dickstein, J., Bengio, S. (2016). *Density estimation using real NVP*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.1605.08803>
- Falkner, S., Coretti, A., Romano, S., Geissler, P., Dellago, C. (2022). *Conditioning normalizing flows for rare event sampling*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.2207.14530>
- Gibson, L.J., & Kroese, D.P. (2022). Rare-event simulation via neural networks. *Advances in modeling and simulation*. Springer. (In Press)
- Glynn, P.W., & Iglehart, D.L. (1989). Importance sampling for stochastic simulations. *Management Science*, 35(11), 1367–1392, <https://doi.org/10.1287/mnsc.35.11.1367>
- Juneja, S., & Shahabuddin, P. (2006). Chapter 11 rare-event simulation techniques: An introduction and recent advances. S.G. Henderson & B.L. Nelson (Eds.), *Simulation* (pp. 291–350). Elsevier.
- Kingma, D.P., & Ba, J. (2017). *Adam: A method for stochastic optimization*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.1412.6980>
- Kingma, D.P., & Dhariwal, P. (2018). Glow: Generative flow with invertible 1x1 convolutions. S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 31). Curran Associates.
- Kobyzev, I., Prince, S.J., Brubaker, M.A. (2021). Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11), 3964–3979, <https://doi.org/10.1109/TPAMI.2020.2992934>

- Kochenderfer, M.J. (2015). *Decision Making Under Uncertainty: Theory and Application*. MIT Press.
- Kroese, D.P., Botev, Z.I., Taimre, T., Vaisman, R. (2019). *Data science and machine learning: Mathematical and statistical methods*. Chapman and Hall/CRC.
- Kroese, D.P., Taimre, T., Botev, Z.I. (2011). *Handbook of Monte Carlo methods*. John Wiley & Sons.
- Liu, J.S. (2004). *Monte Carlo strategies in scientific computing*. Springer.
- Neal, R.M. (2001). Annealed importance sampling. *Statistics and Computing*, 11(2), 125-139, <https://doi.org/10.1023/A:1008923215028>
- Papamakarios, G., Nalisnick, E.T., Rezende, D.J., Mohamed, S., Lakshminarayanan, B. (2021). Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(57), 1-64, Retrieved from <http://jmlr.org/papers/v22/19-1028.html>
- Rezende, D., & Mohamed, S. (2015). Variational inference with normalizing flows. F. Bach & D. Blei (Eds.), *Proceedings of the 32nd international conference on machine learning* (Vol. 37, pp. 1530-1538). PMLR.
- Rubino, G., & Tuffin, B. (2009). *Rare event simulation using Monte Carlo methods*. John Wiley & Sons.
- Rubinstein, R.Y., & Kroese, D.P. (2017). *Simulation and the Monte Carlo method*. John Wiley & Sons.
- Teshima, T., Ishikawa, I., Tojo, K., Oono, K., Ikeda, M., Sugiyama, M. (2020). *Coupling-based invertible neural networks are universal diffeomorphism approximators*. arXiv. Retrieved from <https://doi.org/10.48550/arXiv.2006.11469>
- Tokdar, S.T., & Kass, R.E. (2010). Importance sampling: A review. *WIREs Computational Statistics*, 2(1), 54-60, <https://doi.org/10.1002/wics.56>
- Ziegler, Z., & Rush, A. (2019). Latent normalizing flows for discrete sequences. K. Chaudhuri & R. Salakhutdinov (Eds.), *Proceedings of the 36th international conference on machine learning* (Vol. 97, pp. 7673-7682). PMLR.