

Question 1 Part 1 Report

Yan Kin Chi
yan_kin.chi@unsw.edu.au

I. INTRODUCTION

Efficient Bayesian Inference with Latent Hamiltonian Neural Networks in No-U-Turn Sampling is a paper that builds upon traditional Hamiltonian neural networks (HNNs) but introduces the following novel contributions:

1. Introducing latent variable layers within traditional HNNs to enhance the neural network's expressivity and reduce integration errors.
2. Proposed a threshold-based monitoring scheme that monitors the Hamiltonian gradients inferred by the HNN which reverts to manual integration when errors are large.

II. PAPER CRITIQUE

While the paper is well written, there are various aspects that are unclear that are essential to the understanding/reproduction of the presented results. Below we break down the missing/confusing details into three sections corresponding to the three stages of method. They are generation of the L-HNN training data, training the L-HNN, and conducting Bayesian inference using No-U-Turn sampling (NUTS).

A. Generating training data for L-HNN

- **Stochastic Coverage:** Diversity in trajectories has been emphasized to be essential to avoid training L-HNNs on limited regions of the target density. However, when sampling trajectories, $q(0)$ is defined as the final position of the previous trajectory. This does not maximise the training data's stochastic coverage of the target density and may cause the training data to be biased towards the neighbourhood of a specific mode.
- **Notation Inconsistencies:** Defining the amount of training data to be based on "numerical gradient evaluations" without explicitly defining its meaning makes the exact training details ambiguous. This is made worse with the provided code for generating training samples only based on sample count and trajectory length.
- **Distribution Specific Parameters:** During our testing, parameters for generating training samples are extremely varied for different distributions. However, the paper only explores the critical length (trajectory length) for the 3-D degenerate Rosenbrock density. This makes investigation extremely difficult.
-

B. Training the L-HNN

- **Latent Variables λ :** The paper does not establish a clear structural relationship between the latent variables and the underlying Hamiltonian system. Furthermore, it lacks an explanation of any novel

methodology for integrating latent variables into a standard multilayer perceptron framework. As a result, the proposed L-HNN effectively functions as a HNN with an output layer of λ -dimensions.

Within the paper, there is also no ablation study investigating the effect of modifying λ 's dimensionality on performance. Therefore, the level of contribution of the

- **1-D Gaussian mixture density typo:**

Eq. 18 (1-D Gaussian mixture density) is written as:

$$f(q) \propto 0.5 \exp\left(\frac{(q-1)^2}{2 \times 0.35^2}\right) + 0.5 \exp\left(\frac{(q+1)^2}{2 \times 0.35^2}\right)$$

and is missing a negative sign within each exponent to match the Gaussian density form. Corrected:

$$f(q) \propto 0.5 \exp\left(-\frac{(q-1)^2}{2 \times 0.35^2}\right) + 0.5 \exp\left(-\frac{(q+1)^2}{2 \times 0.35^2}\right)$$

C. Conducting Bayesian Inference using NUTS

- **Excessive use of deterministic parameters:** While the inclusion of algorithm 4 and 5 provides a clear structural representation of the recursive build tree process for NUTS, the use of too many deterministic parameters makes it challenging to evaluate the generalizability of the approach. For example, thresholds Δ_{max}^{hnn} and Δ_{max}^{lf} being set to 10 and 1000 definitely may trigger unnecessary fallbacks (threshold too strict) or fail to detect errors in dynamic regions (too lenient).

III. CODE IMPLEMENTATION

In this section we present investigations results using a code implementation that utilises TensorFlow.

A. *Methodology and Steps to Ensure Correctness:* To minimize implementation errors given the system's complexity, we closely replicated the original code's structure and logic using the TensorFlow framework. While alternative approaches such as leveraging TensorFlow Probability's HMC method exist and warrant future exploration, in this case we focus on maintaining fidelity to the original design.

B. Sampling Experiments

The method's sampling capabilities are mainly demonstrated by sampling the 1-D Gaussian mixture distribution. When training we trained 100 trajectory samples each with $T = 50$ and $\Delta t = 0.05$. To ensure the sample count sufficiently covers all regions of the 1-D Gaussian density, we plotted the position and momentum values for all training samples as shown from in Figure 1.

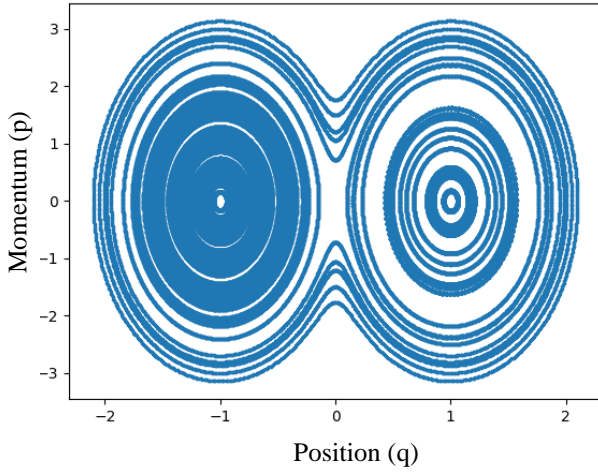


Figure 1: Phase Plot for all training samples

It can be seen that all regions of the density have been adequately covered. For training of the L-HNN we set the hidden dimension to 100, number of hidden layers to be 5, batch size of 1000 and a learning rate of 0.0005.

Next, we generated 8000 samples from the posterior distribution using traditional HMC and L-HNN respectively and plotted the histograms for comparison.

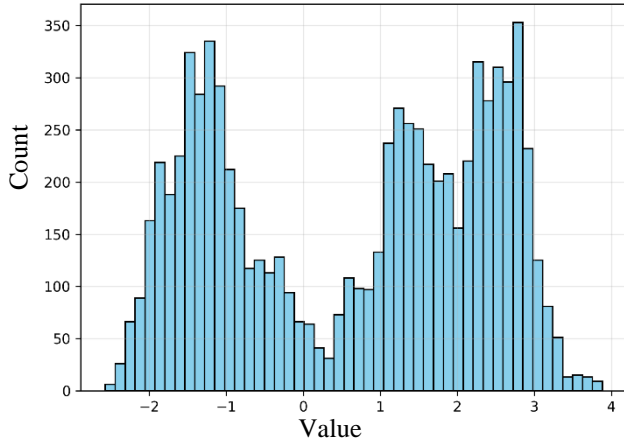


Figure 2: Sampled Distribution (L-HNN)

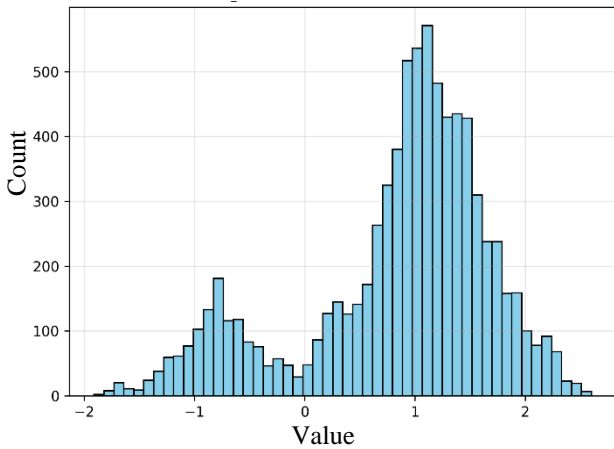


Figure 3: Sampled Distribution (HMC)

For this experiment section we were unable to fully replicate the results shown in Figures 2a, 2b, and 2c of the original paper which show the histograms of HMC and L-HNN closely matching. Specifically, there are a few notable observations:

- **Mode Sample Imbalance:** During both the generation of training samples for L-HNN (Figure 1) and the execution of NUTS, one mode consistently dominates the sample count. This is evident in Figure 1, where the left mode has disproportionately more samples, and in Figures 2 and 3, where the right mode exhibits a similar imbalance. This behaviour is expected because NUTS generates new samples by inheriting the previous position (q) while drawing a new momentum (p) from a standard Gaussian distribution. Transitioning between modes often requires traversing regions of low posterior density, effectively "climbing" over energy barriers, which makes such transitions less likely.

To address this issue and improve exploration, multiple chains can be initiated from each available mode. This approach enhances sampling diversity and ensures better representation of the posterior distribution.

- **Compounding Trajectory Errors:** From our testing, another prominent observation is L-HNNs and by extension HNNs do not perform well when simulating long trajectories.

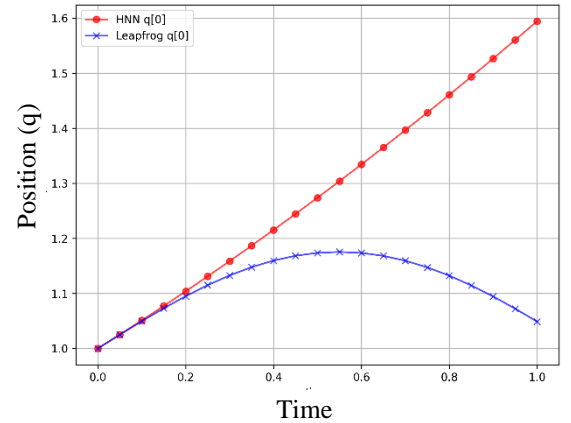


Figure 4: Simulated Trajectory

We demonstrate this by using the same trained model for the previous experiments and start at a random point in phase space. Then we employ a symplectic integrator (i.e. leapfrog) to simulate the trajectory over time using entirely either L-HNNs or HMC. As shown in Figure 4 as time passes the simulated trajectory for L-HNN's deviates further away from HMC. Beyond training optimizations, this suggests that small errors in $\partial H/\partial q$ compound over time.

Overall, this reflects that neural network approaches in Bayesian inference are better suited for tasks requiring a larger number of samples with shorter trajectories. Additionally, the online error monitoring procedure is critical and could be significantly improved by making the threshold more adaptive to the target distribution. Potential improvements setting the threshold dynamically using a moving average or exponential smoothing of Hamiltonian errors from previous samples.