

Bayesian neural network estimation of next-to-leading-order cross sections

by

René Alexander Ask

THESIS

for the degree of

MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences
University of Oslo

Autumn 2021

Bayesian neural network estimation of next-to-leading-order cross sections

René Alexander Ask

© 2021 René Alexander Ask

Bayesian neural network estimation of next-to-leading-order cross sections

<http://www.duo.uio.no/>

Printed: Reprosentralen, University of Oslo

Abstract

This is my abstract.

Acknowledgments

Acknowledgments yo

Contents

Introduction	1
1 Bayesian Learning for Neural Networks	3
1.1 Neural Networks	3
1.1.1 Basic mathematical structure	3
1.1.2 Backpropagation	3
1.1.3 Cost function for regression	4
1.2 Interpretations of probability	5
1.3 Bayesian inference	5
1.4 Bayesian framework for neural networks	5
2 Hamiltonian Monte Carlo	7
2.1 Markov Chain Monte Carlo	7
2.2 Gibbs sampling	8
2.3 The Metropolis-Hastings algorithm	8
2.4 Hamiltonian dynamics	8
2.4.1 Lagrangian Mechanics	9
2.4.2 Hamiltonian Mechanics	9
2.4.3 Leapfrog integration	9
2.5 Hamiltonian Monte Carlo	10
Conclusion	12
Appendices	13
Appendix A	15
A.1 Appendix 1 title	15

Introduction

Motivation, context and problem.

Outline of the Thesis

Give outline of thesis

Chapter 1

Bayesian Learning for Neural Networks

1.1 Neural Networks

In this chapter, I'll introduce the mathematical formalism underpinning neural networks. A neural network is most generally defined as a non-linear function $f : \mathbb{R}^p \rightarrow \mathbb{R}^d$. This non-linear function is built-up as follows:

- An input layer where a real observable $x \in \mathbb{R}^p$ is fed.
- A set of hidden layers.
- An output layer which consists of the output of the function.

1.1.1 Basic mathematical structure

Denote a layer of the network as ℓ . Assume there are L layers in total, such that $\ell \in \{1, 2, \dots, L\}$. Denote the activation at layer ℓ as a_j^ℓ and its bias as b_j^ℓ . Moreover, denote its logits as z_j^ℓ . The following recursive equation defines the general structure of a neural network

$$a_j^\ell = \sigma_\ell \left(\sum_k W_{jk}^\ell a_k^{\ell-1} + b_j^\ell \right) \equiv \sigma_\ell(z_j^\ell), \quad (1.1)$$

where σ_ℓ denotes some possibly non-linear function associated with layer ℓ , W_{jk}^ℓ denotes the weights connecting layer $\ell - 1$ to layer ℓ , and b_j^ℓ denotes the bias at layer ℓ . Typically σ_ℓ is called an activation function. The weights and biases of the network serve as the adjustable or *learnable* parameters of the model.

1.1.2 Backpropagation

The standard approach to train a neural network is minimization of some cost function by employing the *backpropagation* algorithm[1]. The algorithm boils down to four equations defining a recursive algorithm that approximates the gradient with respect to the parameters of the model.

The first of the four equations quantifies the error,

$$\Delta_j^L = \frac{\partial E}{\partial z_j^L}. \quad (1.2)$$

The second equation allows us to compute the error at layer ℓ given we know the error at layer $\ell + 1$,

$$\Delta_j^\ell = \left(\sum_k \Delta_k^{\ell+1} W_{kj}^{\ell+1} \right) \sigma'_\ell(z_j^\ell). \quad (1.3)$$

The final two equations relate these errors to the gradient of the cost function with respect to the model parameters. For the weights, we have

$$\frac{\partial E}{\partial W_{jk}^\ell} = \frac{\partial E}{\partial z_j^\ell} \frac{\partial z_j^\ell}{\partial W_{jk}^\ell} = \Delta_j^\ell a_k^{\ell-1}. \quad (1.4)$$

For the biases, the gradients are

$$\frac{\partial E}{\partial b_j^\ell} = \frac{\partial E}{\partial z_j^\ell} \frac{\partial z_j^\ell}{\partial b_j^\ell} = \Delta_j^\ell. \quad (1.5)$$

With these four equations, we can fit the neural network using simple minimization techniques such as stochastic gradient descent or more complex methods such as ADAM (pages 13-19 in [2]).

1.1.3 Cost function for regression

In this thesis, we're concerned with regression tasks. The activation function of the final layer σ_L is then just the identity function. The typical cost function chosen to solve regression tasks is the L_2 -norm, which for a single output can be written as

$$E(y, \hat{y}) = \frac{1}{2} \|y - \hat{y}\|_2^2, \quad (1.6)$$

where \hat{y} denotes the model output and y the ground-truth. Now, the model output in this case is $\hat{y}_j = a_j^L = z_j^L$. Therefore,

$$\Delta_j^L = \frac{\partial E}{\partial z_j^L} = a_j^L - y_j. \quad (1.7)$$

The backpropagation algorithm for a single datapoint can now be written down. The forward pass is found in algorithm 1.

Algorithm 1 Backpropagation: forward pass

```

 $a_j^0 = x_j$    for  $j = 1, \dots, p$                                  $\triangleright$  Initialize input
for  $l = 1, 2, \dots, L - 1$  do
  for  $j = 1, 2, \dots, n$  do
     $a_j^l \leftarrow \sigma \left( \sum_k W_{jk}^l a_k^{l-1} + b_j^l \right)$ 
  end for
end for
for  $j = 1, 2, \dots, m$  do                                 $\triangleright m$  outputs
   $a_j^L \leftarrow \sigma \left( \sum_k W_{jk}^L a_k^{L-1} + b_j^L \right)$ 
end for

```

The backward pass of the algorithm is stated in algorithm 2.

Algorithm 2 Backpropagation: backward pass

```

for  $j = 1, 2, \dots, m$  do  $\triangleright m$  outputs
     $\Delta_j^L \leftarrow a_j^L - y_j$ 
     $\partial E / \partial b_j^L \leftarrow \Delta_j^L$ 
     $\partial E / \partial W_{jk}^L \leftarrow \Delta_j^L a_k^{L-1}$ 
end for
for  $l = L - 1, \dots, 1$  do
    for  $j = 1, \dots, n$  do
         $\Delta_j^l \leftarrow \left( \sum_k \Delta_k^{l+1} W_{kj}^{l+1} \right) \sigma'(z_j^l)$ 
         $\partial E / \partial b_j^l \leftarrow \Delta_j^l$ 
         $\partial E / \partial W_{jk}^l \leftarrow \Delta_j^l a_k^{l-1}$ 
        update  $b_j^l$  and  $W_{jk}^l$ .
    end for
end for

```

1.2 Interpretations of probability

1.3 Bayesian inference

Given a series of observations $D = \{x^{(1)}, \dots, x^{(n)}\}$, one defines a probabilistic model $P(x^{(i)}, \theta)$ for a set of parameters θ . The foundation for Bayesian inference is Bayes' theorem, which can be formulated as

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \quad (1.8)$$

Some useful terminology is in order. $P(\theta)$ is called the *prior* distribution and embodies our prior knowledge of the distribution of θ before any new observations are considered. $P(D|\theta)$ is called the *likelihood* function and provides information about θ learned from observing the data D . The likelihood function is often written

$$L(\theta) \equiv L(\theta|D) = P(D|\theta). \quad (1.9)$$

One should remember that this is specific to the data D , even though it's sometimes written as $L(\theta)$. The *posterior* distribution $P(\theta|D)$ models our belief about the distribution of θ after observer D . More succinctly, we can write Bayes' theorem as

$$P(\theta|D) \propto L(\theta|D)P(\theta), \quad (1.10)$$

because its rarely of interest, or tractable, to compute $P(D)$, known as the *evidence*.

The objective of Bayesian inference is to compute the *predictive* distribution, which can be expressed as the integral over all parameters of the posterior distribution weighted by the likelihood. Mathematically,

$$P(x^{(n+1)}|D) = \int P(x^{(n+1)}|\theta)P(\theta|D)d\theta = \int L(\theta|x^{(n+1)})P(\theta|D)d\theta. \quad (1.11)$$

1.4 Bayesian framework for neural networks

We can specialize the equations used in Bayesian inference for neural networks in the context of regression. The predictive distribution $P(y|x)$ seeks to model a function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ that for a given x produces an output y . In the infinite data limit, the distribution should be a Dirac Delta function. For finite datasets, however, we instead seek a distribution of outputs given the input features.

Consider a set of observations $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$, where $x^{(i)} \in \mathbb{R}^p$ are the input features and $y^{(i)} \in \mathbb{R}$ are the *targets*. The equation for the predictive distribution changes to

$$P(y^{(n+1)}|x^{(n+1)}, D) = \int L(\theta, y^{(n+1)}|x^{(n+1)})P(\theta|D)d\theta. \quad (1.12)$$

Assuming that the observations in D are drawn independently, the likelihood function can be expressed as

$$L(\theta|D) = \prod_{i=1}^n L(\theta, y^{(i)}|x^{(i)}) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}, \theta). \quad (1.13)$$

In the context of regression, the conditional distribution is chosen to be

$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(f(x) - y)^2}{2\sigma^2}\right), \quad (1.14)$$

where $f(x)$, the neural network output, is the mean of y and σ is the noise of y .

Chapter 2

Hamiltonian Monte Carlo

In this chapter, I'll explain the details of the Hamiltonian Monte Carlo (HMC) method. It's a Monte Carlo (MC) sampling technique that merges Gibbs sampling and a modified version of Metropolis sampling with Hamiltonian dynamics. It avoids the random walk behaviour of the Metropolis algorithm and generates successive samples with smaller correlation. In this chapter, the first thing I'll discuss is Markov chains. From that, follows an outline of Gibbs- and Metropolis sampling. After that, I'll deal with Hamiltonian dynamics before I bring it all together to form the Hamiltonian Monte Carlo algorithm for sampling. This forms the basis for the practical Bayesian learning of the neural networks studied in this thesis.

2.1 Markov Chain Monte Carlo

Monte Carlo Markov chain (MCMC) is a scheme to sample points θ proportional to a distribution $\pi(\theta)$. It generates a new point θ_i given a point θ_{i-1} . By Markov chain, we mean a sequence of points $\theta_1, \theta_2, \dots$ that are possibly dependent, but are each occurring in the sequence in proportion to $\pi(\theta)$. Note that $\pi(\theta)$ here is not an exact probability distribution because it need not be normalized to unity. However, suppose $P(\theta)$ is the underlying probability distribution, then $\pi(\theta) \propto P(\theta)$. Typically, in Bayesian applications, we have a prior $P(\theta)$ and a likelihood $P(D|\theta)$. In this case $\pi(\theta) = P(D|\theta)P(\theta)$ and $\pi(\theta) \propto P(\theta|D)$, that is, it's proportional to the posterior distribution. A few important properties of the Markov chain, originally introduced by Metropolis et. al and built upon by Hastings [3]:

1. **Ergodicity.** Each point θ_i is chosen from a distribution that only depends on the previous point in the sequence, θ_{i-1} . For this, we introduce a transition probability $T(\theta_i|\theta_{i-1})$. Conceptually, means that any point θ can be reached given a long enough sequence of samples.
2. **Detailed balance.** The transition probability is chosen to obey

$$\pi(\theta)T(\theta'|\theta) = \pi(\theta')T(\theta|\theta'),$$

which ensures that the Markov chain is ergodic. Mathematically, we can express this condition as

$$\pi(\theta') = \int \pi(\theta)T(\theta'|\theta)d\theta.$$

3. We allow the transition $\theta \rightarrow \theta$, hence the transition probability $T(\theta|\theta)$ may be non-zero.
4. The transition probabilities integrate to unity, thus

$$\int T(\theta'|\theta)d\theta = 1.$$

5. Finally, the transition probabilities are required to be time-independent.

2.2 Gibbs sampling

Gibbs sampling [4] is a sampling technique for multi-dimensional parameters $\gamma \in \mathbb{R}^d$, for $d > 1$. Suppose $\gamma^{(t)}$ is the parameters at iteration t . Then the parameters $\gamma^{(t+1)}$ at iteration $t + 1$ are generated from $\gamma^{(t)}$ by the algorithm

Algorithm 3 Gibbs sampling

```

Sample  $\gamma_1^{(t+1)} \sim P(\gamma_1 | \gamma_2^{(t)}, \dots, \gamma_d^{(t)})$ 

Sample  $\gamma_2^{(t+1)} \sim P(\gamma_2 | \gamma_1^{(t+1)}, \dots, \gamma_d^{(t)})$ 

 $\vdots$ 

Sample  $\gamma_d^{(t+1)} \sim P(\gamma_d | \gamma_1^{(t+1)}, \dots, \gamma_{d-1}^{(t+1)})$ 

```

2.3 The Metropolis-Hastings algorithm

The Metropolis algorithm [3] is a sampling algorithm based on random walks in parameter space. Albeit efficient for some problems, it's not a suitable sampling technique in the context of neural networks. However, a rudimentary understanding of the algorithm will be useful before we embark upon the HMC sampling algorithm.

The transition probability in the Metropolis algorithm is chosen to be

$$T(\theta' | \theta) = q(\theta' | \theta) A(\theta, \theta'), \quad (2.1)$$

where $q(\theta' | \theta)$ is called the proposal distribution and $A(\theta, \theta')$ is the acceptance probability given by

$$A(\theta, \theta') = \min \left(1, \frac{\pi(\theta') q(\theta | \theta')}{\pi(\theta) q(\theta' | \theta)} \right). \quad (2.2)$$

The point θ' is accepted with probability $A(\theta, \theta')$.

Algorithm 4 Metropolis-Hastings

```

Sample  $\theta' \sim q(\theta' | \theta)$ 

 $p = \min \left( 1, \frac{\pi(\theta') q(\theta | \theta')}{\pi(\theta) q(\theta' | \theta)} \right)$ 

if  $p \geq u$  then
     $\theta \rightarrow \theta'$  ▷ Accept transition
else
     $\theta \rightarrow \theta$  ▷ Reject transition
end if

```

2.4 Hamiltonian dynamics

Hamiltonian dynamics [5] plays a central part in the HMC algorithm. For completeness, I'll first give a brief survey of Lagrangian mechanics from which we derive the Hamiltonian. The Hamiltonian then lays the foundation for the Hamiltonian dynamics.

2.4.1 Lagrangian Mechanics

Assume a set of *generalized coordinates* $q = (q_1, \dots, q_n)$. Generally, the Lagrangian can be written as

$$L(q, \dot{q}, t) = K(q, \dot{q}, t) - V(q, \dot{q}, t), \quad (2.3)$$

where K is the kinetic energy and V is the potential energy of the system. I'll restrict the treatment to the case where there's no explicit dependence on time t . The solutions $q(t)$ can be found by solving the Euler-Lagrange equations given by

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = 0. \quad (2.4)$$

2.4.2 Hamiltonian Mechanics

The Hamiltonian is constructed by the Legendre transformation,

$$H(q, p, t) = \sum_i p_i \dot{q}_i(p_i) - L(q, \dot{q}(p), t), \quad (2.5)$$

where

$$p_i = \frac{\partial L}{\partial \dot{q}_i} \quad (2.6)$$

The equations of motion, known as *Hamilton's equations*, are given by

$$\frac{dp_i}{dt} = \frac{\partial H}{\partial p_i}, \quad \frac{dq_i}{dt} = \frac{\partial H}{\partial q_i}. \quad (2.7)$$

For the purpose of utilizing this framework in the context of HMC, it's assumed that the form of the Lagrangian is

$$L(q, \dot{q}) = K(\dot{q}) - V(q), \quad (2.8)$$

where

$$K(\dot{q}) = \sum_i \frac{1}{2} m_i \dot{q}_i^2. \quad (2.9)$$

The generalized momentum is thus

$$p_i = \frac{\partial K}{\partial \dot{q}_i} = m_i \dot{q}_i, \quad (2.10)$$

from it follows that

$$H(q, p) = K(p) + V(q) = \sum_i \frac{p_i^2}{2m_i} + V(q). \quad (2.11)$$

2.4.3 Leapfrog integration

To run one step of HMC, we need to simulate a Hamiltonian system of the form discussed in the former section. The common choice of algorithm to integrate the equations is *leapfrog* integration [6]. This integrator is *symplectic*, which means it conserves local volumes in phase space. This effectively translates to an approximately conserved value of $H(q, p)$ throughout a simulation, with slight oscillations about a mean value.

Assume we approximate the true coordinates and momenta by (\hat{q}, \hat{p}) . A single leapfrog integration step can then be written as in algorithm 5. h represents the stepsize used in the algorithm.

Algorithm 5 Leapfrog integration (single step)

1. $\hat{p}_i(t + h/2) = \hat{p}_i(t) - \frac{h}{2} \partial V(\hat{q}(t))/\partial q_i$
 2. $\hat{q}_i(t + h) = \hat{q}_i(t) + \frac{h}{m_i} \hat{p}_i(t + h/2)$
 3. $\hat{p}_i(t + h) = \hat{p}_i(t + h/2) - \frac{h}{2} \partial V(\hat{q}(t + h))/\partial q_i$
-

2.5 Hamiltonian Monte Carlo

Hamiltonian Monte Carlo [7] is largely developed and expanded upon by R. Neal. The probability distribution to sample from is expressed in terms the Canonical distribution

$$P(q) \propto \exp(-V(q)), \quad (2.12)$$

where q represents the parameters of the model. $V(q)$ can then always be expressed as

$$V(q) = -\log P(q) - \log Z, \quad (2.13)$$

for some Z . To utilize the framework of Hamiltonian dynamics explained in section 2.4, we introduce momentum variables p_i such that we can expressed the total Hamiltonian as

$$H(q, p) = K(p) + V(q), \quad (2.14)$$

with a corresponding canonical distribution over phase-space

$$P(q, p) \propto \exp(-H(q, p)). \quad (2.15)$$

The algorithm is summarized in

Algorithm 6 Hamiltonian Monte Carlo

```

Randomly choose  $\lambda = \pm 1$ 
Randomly sample  $p$  with Gibbs sampling
Start from current state  $(q, p)$ 
for  $l = 1, \dots, L$  do
     $(q', p') \leftarrow \text{Leapfrog}(q, p)$ 
end for
 $p = \min(1, \exp[-(H(q', p') - H(q, p))])$ 
if  $p \geq u$  then
     $(q, p) \leftarrow (q', p')$ 
else
     $(q, p) \leftarrow (q, p)$ 
end if

```

Conclusion

Conclusion here.

Appendices

Appendix A

A.1 Appendix 1 title

Some appendix stuff.

Bibliography

- [1] D. E. Rumelhart, G. E. Hinton and R. J. Williams, *Learning representations by back-propagating errors*, *Nature* **323** (1986) 533–536.
- [2] P. Mehta, M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher et al., *A high-bias, low-variance introduction to machine learning for physicists*, *Physics Reports* **810** (May, 2019) .
- [3] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller and E. Teller, *Equation of State Calculations by Fast Computing Machines*, *The Journal of Chemical Physics* **21** (1953) 1087–1092, [<https://doi.org/10.1063/1.1699114>].
- [4] S. Geman and D. Geman, *Stochastic relaxation, Gibbs Distributions, and the Bayesian Restoration of Images*, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **PAMI-6** (1984) 721–741.
- [5] J. S. Helbert Goldstein, Charles Poole, *Classical Mechanics*, 3rd ed., ch. 2,8. Addison Wesley, 2000.
- [6] C. M. Bishop, *Pattern Recognition and Machine Learning*, ch. 11, p. 551. Springer New York, 2006.
- [7] R. Neal, *Handbook of Markov Chain Monte Carlo*, ch. 5. Chapman and Hall/CRC, May, 2011. <http://dx.doi.org/10.1201/b10905>. 10.1201/b10905.