
Neural Arithmetic Units

Andreas Madsen^{†‡}
amwebdk@gmail.com

Alexander Rosenberg Johansen[†]
aler@dtu.dk

Who else?

[†]Technical University of Denmark [‡]Computationally Demanding

Abstract

Exact addition, subtraction, multiplication and division present a unique challenge for machine learning models. Neural networks can approximate complex functions by learning from labeled data. However, when extrapolating to out-of-distribution samples on arithmetic operations neural networks often fail to learn the underlying logic, which can be the limiting factor for application such as comparing, counting and inferring physical models. Our proposed Neural Addition Unit (NAU) and Neural Multiplication Unit (NMU) rely on constrained weights to learn rules and extrapolate well beyond the training distribution. The proposed NAU and NMU are inspired by the underlying arithmetic units of the Neural Arithmetic Logic Unit (NALU). The NAU learns addition and subtraction through a linear layer of regularized and constrained weights. The NMU learns multiplication through an accumulative product of the input using gating with an identity function to mask out unwanted elements. Through analytic and empirical analysis we justify how the NAU and NMU improve over the Neural Arithmetic Logic Unit (NALU), a linear regression model and a ReLU based multi-layer perceptron (MLP). Our NAU and NMU have fewer parameters, converges more consistently, learns faster and have more meaningful discrete values than the NALU and its underlying units.

1 Introduction

1

The ability for neurons to hold numbers and do arithmetic operations has been documented in both humans, non-human primates [?], newborn chicks [?] and bees [?]. In our race to solve intelligence we have put much faith in neural networks, which in turn has provided unparalleled and often superhuman performance in many tasks requiring high cognitive ability [???]. However, when using neural networks to solve simple arithmetic problems, such as counting, they systematically fail to extrapolate [???].

In this paper, we analyze and improve parts of the recently proposed Neural Arithmetic Logic Unit (NALU) [?]. Our solution is an alternative formulation of the weight constraint with a clipped linear activation, a regularizer that bias towards sparse solutions, and a reformulation of the multiplica-

¹In the interest of scientific integrity, we have made the code for all experiments, and more, available on GitHub: <https://github.com/AndreasMadsen/stable-nalu>.

tion unit to be partially linear. All of which significantly improves upon the existing NAC_+ and NAC_\bullet units as shown through extensive testing on synthetic tasks and images.

The NALU is a neural network layer with two sub-units. The two sub-units, NAC_+ for addition/subtraction and NAC_\bullet for multiplication/division, are softly gated between with a sigmoid function. By using trainable weights, and restricting the weights towards $\{-1, 0, 1\}$. The weights are learned by observing arithmetic input-output pairs and using backpropagation[?].

We focus only on the NAC_+ and NAC_\bullet as we have found that the gating in NALU can be cumbersome. This is because of the difficulties in selecting between, and simultaneously training, two vastly different operations.

We will thus assume that the appropriate operation is already known, or can empirically be found by varying the network architecture (oracle gating). We find that the NAC_+ and NAC_\bullet units poses optimization difficulties. Our analysis presents the following findings:

- The weight matrix constraint in the NALU, under zero expectation of the mean layer value[?], has a gradient of zero.
- The NAC_\bullet have a treacherous optimization space with unwanted global minimas (as shown in figure 1) and have exploding/vanishing gradients.
- Using the addition module NAC_+ , we observe that the wanted weight matrix values of $\{-1, 0, 1\}$ is rarely found.

Motivated by these convergence and sparsity issue, we propose alternative formulations of the NAC_+ and NAC_\bullet , which we call the Neural Addition Unit (NAU) and Neural Multiplication Unit (NMU).

2 Introducing differentiable binary arithmetic operations

Our goal is to achieve arithmetic operations between the elements of a vector. Formally defined as .

$$x_1 \circ_1 x_2 \circ_2 \cdots x_{k-1} \circ_{k-1} x_k \mid (x_1, \dots, x_k) \in \mathbf{x}, \mathbf{x} \in \mathbb{R}^n, \circ \in \{+, -, \cdot, \div\} \quad (1)$$

The Neural Arithmetic Logic Unit (NALU)[?] attempts to solve equation 1 by presenting two sub-units; the NAC_+ and NAC_\bullet to exclusively represent either the $\{+, -\}$ or the $\{\cdot, \div\}$ operations. The NALU attempts to have either NAC_+ or NAC_\bullet selected exclusively, which could require the NALU to be applied multiple times (alternating between NAC_+ and NAC_\bullet) in order to represent the entire space of solutions for equation 1.

The NAC_+ and NAC_\bullet are defined accordingly,

$$W_{h_\ell, h_{\ell-1}} = \tanh(\hat{W}_{h_\ell, h_{\ell-1}}) \sigma(\hat{M}_{h_\ell, h_{\ell-1}}) \quad (2)$$

$$\text{NAC}_+ : z_{h_\ell} = \sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} z_{h_{\ell-1}} \quad (3)$$

$$\text{NAC}_\bullet : z_{h_\ell} = \exp \left(\sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon) \right) \quad (4)$$

where $\hat{\mathbf{W}}, \hat{\mathbf{M}} \in \mathbb{R}^{H_\ell \times H_{\ell-1}}$ are trainable weight matrices. These are then combined using tanh and sigmoid transformation to bias the parameters towards a $\{-1, 0, 1\}$ solution. Having $\{-1, 0, 1\}$ allows a linear layer to exactly emulate the binary $\{+, -\}$ operation between elements of a vector as used when computing the NAC_+ . The NAC_\bullet extends the NAC_+ by using an exponential log transformation, which, with $\{-1, 0, 1\}$ weight values, becomes the $\{\cdot, \div\}$ operations (within ϵ precision).

The NALU combines these units with a gating mechanism $\mathbf{z} = \mathbf{g} \odot \text{NAC}_+ + (1 - \mathbf{g}) \odot \text{NAC}_\bullet$ given $\mathbf{g} = \sigma(\mathbf{G}\mathbf{x})$. The idea is that the NALU should be a plug-and-play component in a neural network and has the ability to, with stochastic gradient descent and backpropagation, to learn the functionality in equation 1.

We have spent a lot of time on gating, and have an intuition on why it might not work. It turns out to be extremely difficult. Unfortunately we have no solution. How do we argue that we don't look at gating?

Can this be more formally defined

maybe visualize the function space

2.1 Challenges of the NALU, NAC_+ and NAC_\bullet

To simplify the problem we have chosen to leave out the gating mechanism and focus on the sub-units, assuming "oracle gating".

We find this reasonable given ...

We find that gating between NAC_+ and NAC_\bullet is challenging to make converge and it is not obvious whether or not the gate could bias towards a specific sub-unit. We have not had any consistent success of convergence using the gating mechanism using the NALU or by combining our own proposed sub-units (NAU, NMU). We find that gating between NAC_+ and NAC_\bullet is challenging. This is likely due to the vastly different gradients, causing addition to be learned much faster than multiplication.

In section 2.1.1 we analyse the gradients of equation 2 and propose a clipped linear alternative, with a regularizer that biases towards $\{-1, 0, 1\}$. This results in an alternative addition unit called NAU.

In section 2.1.2 we analyse the gradient of NAC_\bullet , and propose an alternative construction called NMU, that has a more well-behaved loss space.

In section 2.1.4 we analyse the moments of NAC_+ and NAU, and derive the optimal weight initializations.

In section 2.1.5 we analyse the moments of NAC_\bullet and NMU, and show that the NMU have a more well-behaved expectation and variance.

Maybe remove this overview?

2.1.1 Weight matrix construction

The weight matrix construction $\tanh(\hat{W}_{h_{\ell-1}, h_\ell})\sigma(\hat{M}_{h_{\ell-1}, h_\ell})$ has the following properties that could make convergence challenging using gradient decent.

The loss gradient with respect to the weight matrices can be derived from equation 2.

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \hat{W}_{h_{\ell-1}, h_\ell}} &= \frac{\partial \mathcal{L}}{\partial W_{h_{\ell-1}, h_\ell}} (1 - \tanh^2(\hat{W}_{h_{\ell-1}, h_\ell})) \sigma(\hat{M}_{h_{\ell-1}, h_\ell}) \\ \frac{\partial \mathcal{L}}{\partial \hat{M}_{h_{\ell-1}, h_\ell}} &= \frac{\partial \mathcal{L}}{\partial W_{h_{\ell-1}, h_\ell}} \tanh(\hat{W}_{h_{\ell-1}, h_\ell}) \sigma(\hat{M}_{h_{\ell-1}, h_\ell}) (1 - \sigma(\hat{M}_{h_{\ell-1}, h_\ell}))\end{aligned}\quad (5)$$

The gradient $E[\partial \mathcal{L} / \partial \hat{M}_{h_{\ell-1}, h_\ell}] = 0$ can be problematic as we prefer zero having a zero mean expectation of our output. Something that can only be ensured with $E[\hat{W}_{h_{\ell-1}, h_\ell}] = 0$ [?].

In our empirical analysis we find that equation 2 does not create the desired bias for $\{-1, 0, 1\}$, as it doesn't converge towards those values.

To create a bias and prevent the gradient challenges of equation 5 we propose a simple clamped linear construction with an out-of-bound regularizer $\mathcal{R}_{\ell, \text{oob}}$ to force \hat{W} to be within $[-1, 1]$ and ensure that the gradient is always present.

$$\begin{aligned}W_{h_{\ell-1}, h_\ell} &= \min(\max(\hat{W}_{h_{\ell-1}, h_\ell}, -1), 1), \\ \mathcal{R}_{\ell, \text{bias}} &= \frac{1}{H_\ell + H_{\ell-1}} \sum_{h_{\ell-1}=1}^{H_\ell} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \hat{W}_{h_{\ell-1}, h_\ell}^2 (1 - |\hat{W}_{h_{\ell-1}, h_\ell}|)^2 \\ \mathcal{R}_{\ell, \text{oob}} &= \frac{1}{H_\ell + H_{\ell-1}} \sum_{h_{\ell-1}=1}^{H_\ell} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \max(|\hat{W}_{h_{\ell-1}, h_\ell}| - 1, 0)^2 \\ \text{NAU : } z_{h_\ell} &= \sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} z_{h_{\ell-1}} \\ \mathcal{L} &= \hat{\mathcal{L}} + \lambda_{\text{bias}} \mathcal{R}_{\ell, \text{bias}} + \lambda_{\text{oob}} \mathcal{R}_{\ell, \text{oob}}\end{aligned}\quad (6)$$

2.1.2 Challenges of division

The NAC_\bullet , as formulated in equation 4, has the ability to learn exact multiplication and division of elements from a vector if the weights of $W_{h_{\ell-1}, h_\ell}$ are one of $\{-1, 0, 1\}$.

However, backpropagation through the NAC_\bullet unit reveals that if $|z_{h_{\ell-1}}|$ is near zero, $W_{h_{\ell-1},h_\ell}$ is negative and ϵ is small, the gradient term will explode and oscillate between large positive and large negative values, which can be problematic in optimization [?], as visualized in figure 1.

$$\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} = \sum_{h_\ell=1}^{H_\ell} \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \frac{\partial z_{h_\ell}}{\partial z_{h_{\ell-1}}} = \sum_{h_\ell=1}^{H_\ell} \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} z_{h_\ell} W_{h_\ell,h_{\ell-1}} \frac{\text{sign}(z_{h_{\ell-1}})}{|z_{h_{\ell-1}}| + \epsilon} \quad (7)$$

(see full derivation in Appendix A.2)

This is not an issue for positive values of $W_{h_{\ell-1},h_\ell}$ (multiplication), as z_{h_ℓ} and $z_{h_{\ell-1}}$ will be correlated causing the terms z_{h_ℓ} and $\frac{\text{sign}(z_{h_{\ell-1}})}{|z_{h_{\ell-1}}| + \epsilon}$ to cancel out.

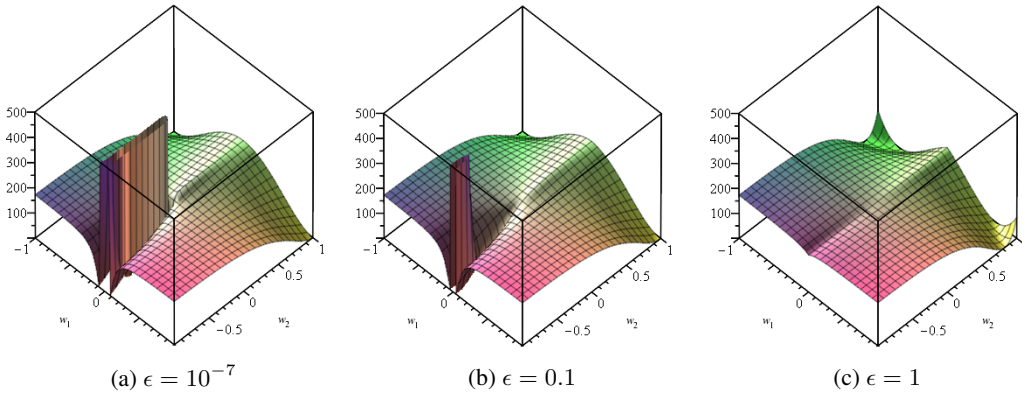


Figure 1: RMS loss curvature for a NAC_+ layer followed by a NAC_\bullet layer. The weight matrices constrained are to $\mathbf{W}_1 = \begin{bmatrix} w_1 & w_1 & 0 & 0 \\ w_1 & w_1 & w_1 & w_1 \end{bmatrix}$, $\mathbf{W}_2 = \begin{bmatrix} w_2 & w_2 \end{bmatrix}$. The problem is $x = (1, 1.2, 1.8, 2)$, $t = 13.2$. Desired solution is $w_1 = w_2 = 1$, although this problem have additional undesired solutions.

These observations are particular problematic when considering that $E[z_{h_{\ell-1}}] = 0$ is a desired property when initializing [?]. An alternative multiplication operator must thus be able to not explode for $z_{h_{\ell-1}}$ near zero. To that end we propose a new neural multiplication units (NMU):

$$\begin{aligned} W_{h_{\ell-1},h_\ell} &= \min(\max(\hat{W}_{h_{\ell-1},h_\ell}, 0), 1), \\ \mathcal{R}_{\ell,\text{bias}} &= \frac{1}{H_\ell + H_{\ell-1}} \sum_{h_\ell=1}^{H_\ell} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \hat{W}_{h_{\ell-1},h_\ell}^2 (1 - \hat{W}_{h_{\ell-1},h_\ell})^2 \\ \mathcal{R}_{\ell,\text{oob}} &= \frac{1}{H_\ell + H_{\ell-1}} \sum_{h_\ell=1}^{H_\ell} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \max\left(\left|\hat{W}_{h_{\ell-1},h_\ell} - \frac{1}{2}\right| - \frac{1}{2}, 0\right)^2 \\ \text{NMU} : z_{h_\ell} &= \prod_{h_{\ell-1}=1}^{H_{\ell-1}} (W_{h_{\ell-1},h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1},h_\ell}) \end{aligned} \quad (8)$$

This units does not support division, but supporting division is likely infeasible if $z_{h_{\ell-1}}$ near zero should not cause explosions. The NALU paper also shows that division doesn't work well for their unit, hence very little is lost here [?]. On the other hand, this unit construction understand the difference between a negative and a positive $z_{h_{\ell-1}}$ values, which should be considered an added bonus, as this allows extrapolations into the negative input range.

The gradients weight gradient and backpropagation term of the NMU are (see details in Appendix A.3):

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}} &= \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \frac{\partial z_{h_\ell}}{\partial W_{h_\ell, h_{\ell-1}}} = \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \frac{z_{h_\ell}}{W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell}} (z_{h_{\ell-1}} - 1) \\ \frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} &= \sum_{h_\ell=1}^{H_\ell} \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \frac{\partial z_{h_\ell}}{\partial z_{h_{\ell-1}}} = \sum_{h_\ell=1}^{H_\ell} \frac{z_{h_\ell}}{W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell}} W_{h_{\ell-1}, h_\ell}\end{aligned}\quad (9)$$

These is much more well-behaved. Note also that the fraction does not explode for $z_{h_{\ell-1}}$ close to zero, as the denominator simply cancels out a term in z_{h_ℓ} .

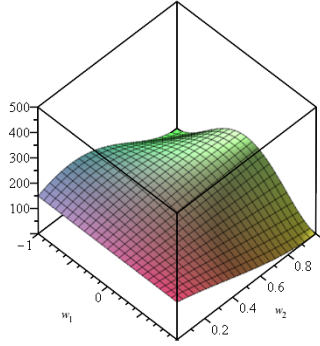


Figure 2: RMS loss curvature (without regularization) for a NAC_+ layer followed by an NMU layer. Otherwise, the setup is identical to that in Figure 1.

2.1.3 Moments and initialization

Initialization is important to consider for fast and consistent convergence [?].

Our proposed NAU, can be initialize using Glorot initialization as it is just a linear layer. The NAC_+ unit can also achieve an ideal initialization, although it is less trivial (details in Appendix B.1).

Using second order multivariate Taylor approximation and some assumptions of uncorrelated stochastic variables, the expectation of NAC_\bullet can be estimated to be

$$E[z_{h_\ell}] \approx \left(1 + \frac{1}{2} \text{Var}[W_{h_\ell, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2\right)^{H_{\ell-1}} \Rightarrow E[z_{h_\ell}] > 1 \quad (10)$$

(proof in Appendix B.2). An ideal initialization should satisfy $E[z_{h_\ell}] = 0$ [?], which the expectation for NAC_\bullet can't.

Our proposed NMU when initialized with $E[W_{h_\ell, h_{\ell-1}}] = 1/2$ has an expectation of

$$E[z_{h_\ell}] \approx \left(\frac{1}{2}\right)^{H_{\ell-1}} \quad (11)$$

which is not zero, but does approach zero for $H_{\ell-1} \rightarrow \infty$ (proof in Appendix B.3).

The variance is also more well-behaved, but does similarly to NAC_\bullet not provide a input-independent initialization strategy (proof in Appendix B.2 and B.3). We propose initializing with $\text{Var}[W_{h_{\ell-1}, h_\ell}] = \frac{1}{4}$, as this provides an ideal initialization if $\text{Var}[z_{h_{\ell-1}}] = 1$ and $H_{\ell-1}$ is large (see proof in Appendix B.3.3).

2.1.4 Moments and initialization for addition

Maybe put
this entire
section in
appendix?

Initialization is important for fast and consistent convergence. The desired properties are according to Glorot et al. [?]:

$$\begin{aligned} E[z_{h_\ell}] &= 0 & E\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &= 0 \\ \text{Var}[z_{h_\ell}] &= \text{Var}[z_{h_{\ell-1}}] & \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &= \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}}\right] \end{aligned} \quad (12)$$

The NAU layer is trivial, as this is just a linear layer. Thus the result from Glorot et al. ($\text{Var}[W_{h_{\ell-1}, h_\ell}] = \frac{2}{H_{\ell-1} + H_\ell}$) can be used [?].

However, the original NAC₊ unit is less trivial as $W_{h_{\ell-1}, h_\ell}$ is not sampled directly. Assuming that $\hat{W}_{h_\ell, h_{\ell-1}} \sim \text{Uniform}[-r, r]$ and $\hat{M}_{h_\ell, h_{\ell-1}} \sim \text{Uniform}[-r, r]$ then the variance can be derived (see proof in Appendix B.1) to be:

$$\text{Var}[W_{h_{\ell-1}, h_\ell}] = \frac{1}{2r} \left(1 - \frac{\tanh(r)}{r}\right) \left(r - \tanh\left(\frac{r}{2}\right)\right) \quad (13)$$

One can solve for r , given the desired variance.

2.1.5 Moments and initialization for multiplication

Using second order multivariate Taylor approximation and some assumptions of uncorrelated stochastic variables, the expectation and variance of the NAC_• layer can be estimated to:

$$\begin{aligned} f(c_1, c_2) &= \left(1 + c_1 \frac{1}{2} \text{Var}[W_{h_\ell, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2\right)^{c_2 H_{\ell-1}} \\ E[z_{h_\ell}] &\approx f(1, 1) \\ \text{Var}[z_{h_2}] &\approx f(4, 1) - f(1, 2) \\ E\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &= 0 \\ \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &\approx \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}}\right] H_\ell f(4, 1) \text{Var}[W_{h_\ell, h_{\ell-1}}] \\ &\quad \cdot \left(\frac{1}{(|E[z_{h_{\ell-1}}]| + \epsilon)^2} + \frac{3}{(|E[z_{h_{\ell-1}}]| + \epsilon)^4} \text{Var}[z_{h_{\ell-1}}]\right) \end{aligned} \quad (14)$$

Maybe put this entire section in appendix?

put in assumptions in appendix

This is problematic because $E[z_{h_\ell}] \geq 1$, and the variance explodes for $E[z_{h_{\ell-1}}] = 0$ which is normally a desired property.

For our proposed NMU, the expectation and variance can be derived (see proof in Appendix B.3) using the same assumptions as before, although no Taylor approximation is required:

$$\begin{aligned} E[z_{h_\ell}] &\approx \left(\frac{1}{2}\right)^{H_{\ell-1}} \\ E\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &\approx 0 \\ \text{Var}[z_{h_\ell}] &\approx \left(\text{Var}[W_{h_{\ell-1}, h_\ell}] + \frac{1}{4}\right)^{H_{\ell-1}} (\text{Var}[z_{h_{\ell-1}}] + 1)^{H_{\ell-1}} - \left(\frac{1}{4}\right)^{H_{\ell-1}} \\ \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}}\right] &\approx \text{Var}\left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}}\right] H_\ell \\ &\quad \cdot \left(\left(\text{Var}[W_{h_{\ell-1}, h_\ell}] + \frac{1}{4}\right)^{H_{\ell-1}} (\text{Var}[z_{h_{\ell-1}}] + 1)^{H_{\ell-1}-1} - \left(\frac{1}{4}\right)^{H_{\ell-1}}\right) \end{aligned} \quad (15)$$

consider
throwing in
appendix

These expectations are much more well behaved. It is properly unlikely to expect that the expectation can become zero, since the identity for multiplication is 1. However, for a large $H_{\ell-1}$ it will be near zero.

The variance is also more well-behaved, but does not provide a input-independent initialization strategy. We propose initializing with $Var[W_{h_{\ell-1}, h_{\ell}}] = \frac{1}{4}$, as this is the solution to $Var[z_{h_{\ell}}] = Var[z_{h_{\ell-1}}]$ assuming $Var[z_{h_{\ell-1}}] = 1$ and a large $H_{\ell-1}$ (see proof in Appendix B.3.3). However, feel free to compute more exact solutions.

3 Experimental results

3.1 Simple function task

3.1.1 Task definition

Our simple function task samples an input vector \mathbf{x} from a uniform distribution. From this input vector, the sum of two subsets a and b are then computed. Finally the target t is then an operation performed on a and b (e.g. $a \cdot b$). This is identical to the task by the same name in the Original NALU paper [?]. Except that we parameterize it in order to compare the models for different configurations, see figure 3. To make comparison simple, we define a set of default parameters (table 1) and only vary one of them at the time.

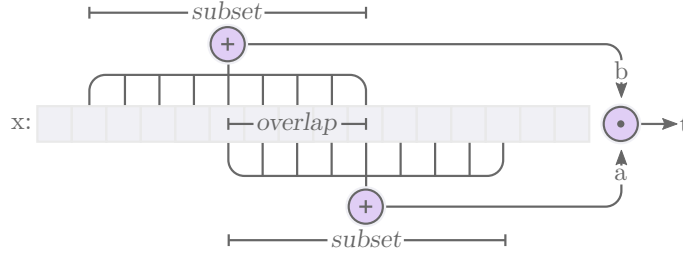


Figure 3: Dataset is parameterized into “Input Size”, “Subset Ratio”, “Overlap Ratio”, an Operation (here showing multiplication), “Interpolation Range” and “Extrapolation Range” from which the data set sampled.

Table 1: Default dataset parameters

| Parameter Name | Default Value |
|---------------------|---------------|
| Input Size | 100 |
| Subset Ratio | 0.25 |
| Overlap Ratio | 0.5 |
| Interpolation Range | $U[1, 2]$ |
| Extrapolation Range | $U[1, 6]$ |

The training set is sampled from the interpolation range throughout training. The extrapolation range is defined as 2048 samples picked for each seed at the start of training.

should ex-
trapolation be
 $U[2, 6]$?

3.1.2 Criterion

Normally one would report the interpolation and extrapolation loss. However, the complex approximations that one would typically see in neural networks are not considered good enough. The goal is to achieve a solution that is sufficiently close to a perfect solution. Because there can be many valid permutations of a perfect solution, especially for addition, a solution is judged first on the final extrapolation error, and then on a sparsity error. Because the interpolation and extrapolation

should we
define valid
and test sets?
maybe we
should make
one BIG test
set and have
all models
run on that

errors are quite noisy, even for a near perfect solution. The median over the last 100 measurements is reported.

A model is considered a success if the extrapolation median is less than $\epsilon = 0.2$. This value was acquired by inspecting the error of a near perfect solution.

The sparsity error is computed as in equation 16, and is only considered for the models that did solve the last.

$$E_{\text{sparsity}} = \max_{h_{\ell-1}, h_{\ell}} \min(|W_{h_{\ell-1}, h_{\ell}}|, |1 - |W_{h_{\ell-1}, h_{\ell}}||) \quad (16)$$

The first iteration for which extrapolation $< \epsilon$, is also reported. Again, only models that did solve the task are considered.

use the 2048 as a validation set and just pick the best instead. Then report test results for those

Get a better criterion.

3.1.3 Model setup

For all experiments the $\mathcal{R}_{\ell, \text{obb}}$ regularizer is added to the loss without modification or scaling, while the $\mathcal{R}_{\ell, \text{bias}}$ regularizer is gradually upscaled with $0.1 \cdot (1 - \exp(-10^5 \cdot t))$. Generally this regularizer should just be sufficiently small to not interfere with early training.

3.1.4 Very simple function

To empirically validate the theoretical problems with NAC_{\bullet} , let's consider the very simple problem shown earlier in figure 1. That is $x \in \mathbb{R}^4$, $a = x_1 + x_2$ and $b = x_1 + x_2 + x_3 + x_4$. The solution to this problem is that seen in equation 17.

$$\mathbf{W}_1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \mathbf{W}_2 = [1 \quad 1] \quad (17)$$

Each model is trained 100 times with different seeds, and stopped after 200000 iterations. Default Adam optimization is used, with a mini-batch size of 128 observations. The results (as seen in table 2), shows that NMU have a much higher success rate and converges much faster. The few cases that did not converge successfully are because of underflow when $w = 0$ in the NMU layer.

Table 2: Shows the success-rate for extrapolation $< \epsilon$, at what global step the model converged at, and the sparsity error for all weight matrices.

| Operation | Model | Success rate | Converged at | Sparsity error |
|-------------|------------------------|--------------|--------------|----------------------|
| $a \cdot b$ | NAC_{\bullet} | 13% | 43923 | 7.5×10^{-6} |
| | NALU | 26% | 54115 | 9.2×10^{-6} |
| | NMU | 94% | 17415 | 1.4×10^{-4} |

3.1.5 Defaults

To compare on the exact same task as used in the Original NALU paper [?]. We report the success rate, the iteration which the model converged, and the sparsity error in table 3. The models are trained for 5000000 iterations. Default Adam optimization is used, with a mini-batch size of 128 observations. The NMU model is an NAU layer followed by an NMU layer. Likewise the NAC_{\bullet} model, is a NAC_{+} layer followed by a NAC_{\bullet} layer.

As seen the NMU model, unlike the NAC_{\bullet} model always converges, and even when NAC_{\bullet} model converges the NMU models converges about twice as fast.

The NAU model, like the NAC_{+} model, always converges. However, NAU model converges more than twice as fast. It even converges faster than a Linear model. Also notice that the NAC_{+} model have a poor sparsity error. This is because it doesn't bias to $\{-1, 0, -1\}$.

Table 3: Shows the success-rate for extrapolation $< \epsilon$, at what global step the model converged at, and the sparsity error for all weight matrices.

| Operation | Model | Success rate | Converged at | Sparsity error |
|-------------|------------------|--------------|--------------|----------------------|
| $a \cdot b$ | NAC $_{\bullet}$ | 30% | 2529667 | 3.9×10^{-4} |
| | Linear | 0% | — | — |
| | NALU | 0% | — | — |
| | NMU | 90% | 1619333 | 1.8×10^{-3} |
| $a - b$ | NAC $_{+}$ | 100% | 1479400 | 4.6×10^{-1} |
| | Linear | 100% | 1854600 | 3.7×10^{-1} |
| | NALU | 20% | 3618000 | 4.7×10^{-1} |
| | NAU | 100% | 4200 | 1.9×10^{-3} |
| $a + b$ | NAC $_{+}$ | 100% | 71200 | 4.8×10^{-1} |
| | Linear | 100% | 42200 | 6.1×10^{-1} |
| | NALU | 0% | — | — |
| | NAU | 100% | 703700 | 1.7×10^{-3} |

add linear to multiplication

3.1.6 Exploration of dataset parameters

Finally, the parameters from which the dataset is constructed are considered for just the multiplication problem ($a \cdot b$). The setup is the same the results from table 3. The results are visualized in figure 4, 5, 7, and 6. Errors bars show the upper and lower 10% quantile, computed over 10 different seeds for each configuration. The center shows the mean of those 10 observations.

Generally the NMU performs far better than both NAC $_{\bullet}$ and NALU. Some important observations to make:

- Input size > 100 . The NMU model’s success-rate very suddenly decreases when the input size is greater than 100. We have been unable to explain why this happens. We suspect it is a problem with the signal-to-noise ratio of the problem. However the result is also seen if the mini-batch size is dramatically increased.
- Overlap ratio = 0: Both the NMU and also the NAC $_{\bullet}$ when it does converge, finds a suboptimal solution where in the addition layer $w = 1$ for the overlapping input between a and b , and $w = 0$ for where the input isn’t used. However when an input-scalar is only used in either a or b , convergence the corresponding weights is difficult and slow. Thus the lower the overlap ratio is, the harder the problem is.

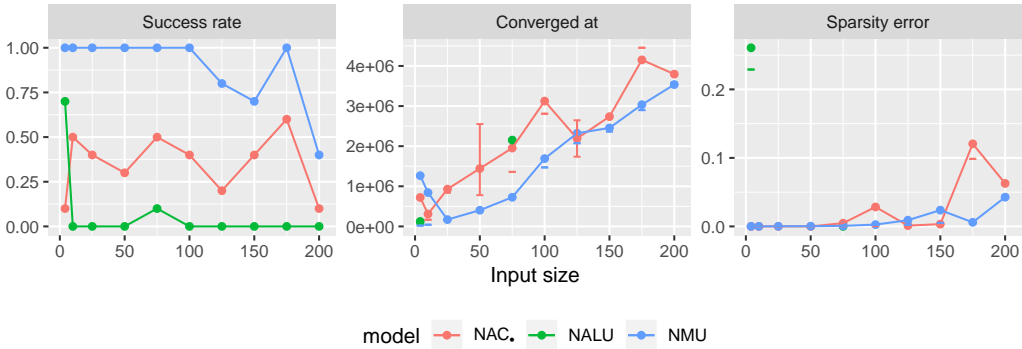


Figure 4: Shows the effect of the input size, on the simple function task problem.

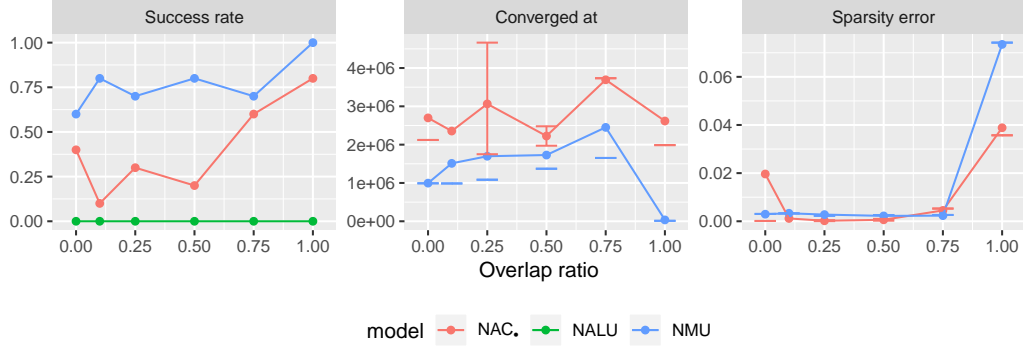


Figure 5: Shows the effect of the overlap ratio, on the simple function task problem.

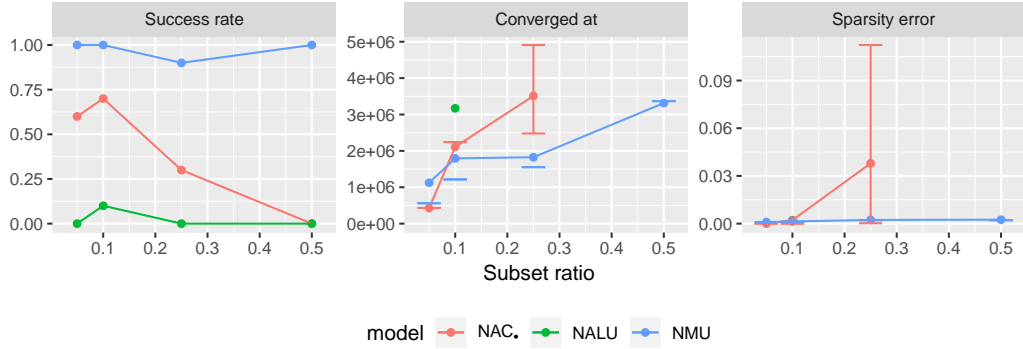


Figure 6: Shows the effect of the subset ratio, on the simple function task problem.

3.2 Sequential MNIST

To evaluate NAU and NMU in a end-to-end context in combination with a more complex network. We consider the Sequential MNIST Arithmetic task, also presented in the NALU paper [?].

The task is to take a sequence of MNIST images, then use a CNN layer² to produce a hidden layer with 10 elements which somehow describes the number. An recurrent arithmetic unit, is then used to either sum or multiply each MNIST digit together.

This is slightly different from [?], as they only considered addition in the form of counting and a sum. While here consider the sum and the product of a sequence is a considered (figure 8). Such that the multiplication layer also can be judged.

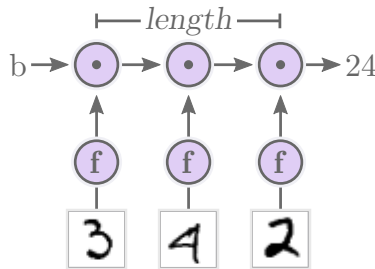


Figure 8: Shows how $3 \cdot 4 \cdot 2$ is computed from a sequence of MNIST digits. f is a CNN layer, that transforms the image into a 10-long vector representation.

²We use the model from <https://github.com/pytorch/examples/tree/master/mnist> which is also used in the NALU paper, [?].

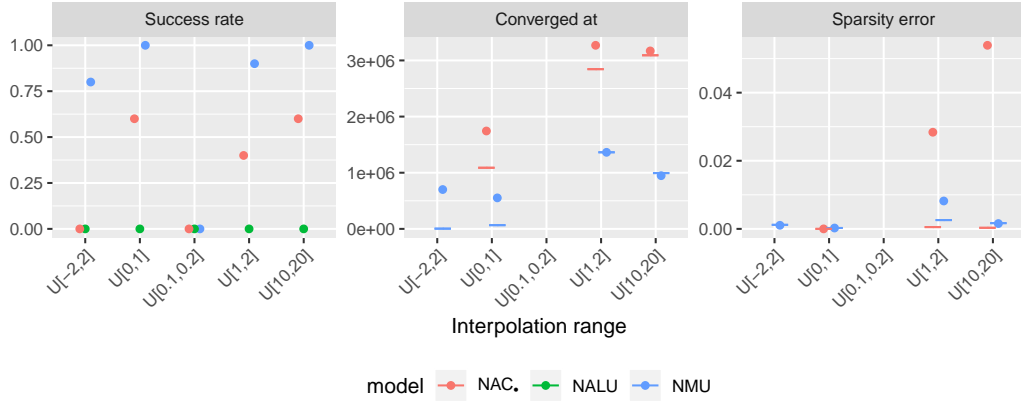


Figure 7: Shows the effect of the interpolation range. For each interpolation range, the following extrapolation ranges are used: $U[-2, 2] \rightarrow U[-6, 6]$, $U[0, 1] \rightarrow U[0, 5]$, $U[0.1, 0.2] \rightarrow U[0, 2]$, $U[1, 2] \rightarrow U[1, 6]$, $U[10, 20] \rightarrow U[1, 40]$.

Still waiting for results to be computed.

4 Future work

- Lorem Ipsum

5 Conclusion

Acknowledgments

We would like to thank Andrew Trask and the other authors of the NALU paper, for highlighting the importance and challenges of etrapolation in Neural Networks. We would also like to thank the students Raja Shan Zaker Kreen and William Frisch Moller from The Technical University of Denmark, who showed us that the NALU does not converge consistently.

A Gradient derivatives

A.1 Weight matrix construction

For clarity the weight matrix construction is defined using scalar notation

$$W_{h_\ell, h_{\ell-1}} = \tanh(\hat{W}_{h_\ell, h_{\ell-1}}) \sigma(\hat{M}_{h_\ell, h_{\ell-1}}) \quad (18)$$

The of the loss with respect to $\hat{W}_{h_\ell, h_{\ell-1}}$ and $\hat{M}_{h_\ell, h_{\ell-1}}$ is then straight forward to derive.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \hat{W}_{h_\ell, h_{\ell-1}}} &= \frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}} \frac{\partial W_{h_\ell, h_{\ell-1}}}{\partial \hat{W}_{h_\ell, h_{\ell-1}}} \\ &= \frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}} (1 - \tanh^2(\hat{W}_{h_\ell, h_{\ell-1}})) \sigma(\hat{M}_{h_\ell, h_{\ell-1}}) \\ \frac{\partial \mathcal{L}}{\partial \hat{M}_{h_\ell, h_{\ell-1}}} &= \frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}} \frac{\partial W_{h_\ell, h_{\ell-1}}}{\partial \hat{M}_{h_\ell, h_{\ell-1}}} \\ &= \frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}} \tanh(\hat{W}_{h_\ell, h_{\ell-1}}) \sigma(\hat{M}_{h_\ell, h_{\ell-1}}) (1 - \sigma(\hat{M}_{h_\ell, h_{\ell-1}})) \end{aligned} \quad (19)$$

As seen from this result, one only needs to consider $\frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}}$ for NAC_+ and NAC_\bullet , as the gradient with respect to $\hat{W}_{h_\ell, h_{\ell-1}}$ and $\hat{M}_{h_\ell, h_{\ell-1}}$ is just a multiplication on $\frac{\partial \mathcal{L}}{\partial W_{h_\ell, h_{\ell-1}}}$.

A.2 Gradient of NAC_\bullet

First the NAC_\bullet is defined using scalar notation.

$$z_{h_\ell} = \exp \left(\sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon) \right) \quad (20)$$

The gradient of the loss with respect to $W_{h_\ell, h_{\ell-1}}$ is straight forward to derive.

$$\begin{aligned} \frac{\partial z_{h_\ell}}{\partial W_{h_\ell, h_{\ell-1}}} &= \exp \left(\sum_{h'_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h'_{\ell-1}} \log(|z_{h'_{\ell-1}}| + \epsilon) \right) \log(|z_{h_{\ell-1}}| + \epsilon) \\ &= z_{h_\ell} \log(|z_{h_{\ell-1}}| + \epsilon) \end{aligned} \quad (21)$$

We now wish to derive the backpropagation term $\delta_{h_\ell} = \frac{\partial \mathcal{L}}{\partial z_{h_\ell}}$, because z_{h_ℓ} affects $\{z_{h_{\ell+1}}\}_{h_{\ell+1}=1}^{H_{\ell+1}}$ this becomes:

$$\delta_{h_\ell} = \frac{\partial \mathcal{L}}{\partial z_{h_\ell}} = \sum_{h_{\ell+1}=1}^{H_{\ell+1}} \frac{\partial \mathcal{L}}{\partial z_{h_{\ell+1}}} \frac{\partial z_{h_{\ell+1}}}{\partial z_{h_\ell}} = \sum_{h_{\ell+1}=1}^{H_{\ell+1}} \delta_{h_{\ell+1}} \frac{\partial z_{h_{\ell+1}}}{\partial z_{h_\ell}} \quad (22)$$

To make it easier to derive $\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_\ell}}$ we re-express the z_{h_ℓ} as $z_{h_{\ell+1}}$.

$$z_{h_{\ell+1}} = \exp \left(\sum_{h_\ell=1}^{H_\ell} W_{h_{\ell+1}, h_\ell} \log(|z_{h_\ell}| + \epsilon) \right) \quad (23)$$

The gradient of $\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}}$ is then:

$$\begin{aligned}\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}} &= \exp\left(\sum_{h_{\ell}=1}^{H_{\ell}} W_{h_{\ell+1}, h_{\ell}} \log(|z_{h_{\ell}}| + \epsilon)\right) W_{h_{\ell+1}, h_{\ell}} \frac{\partial \log(|z_{h_{\ell}}| + \epsilon)}{\partial z_{h_{\ell}}} \\ &= \exp\left(\sum_{h_{\ell}=1}^{H_{\ell}} W_{h_{\ell+1}, h_{\ell}} \log(|z_{h_{\ell}}| + \epsilon)\right) W_{h_{\ell+1}, h_{\ell}} \frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \\ &= m_{h_{\ell+1}} W_{h_{\ell+1}, h_{\ell}} \frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon}\end{aligned}\quad (24)$$

$\text{abs}'(z_{h_{\ell}})$ is the gradient of the absolute function. In the paper we denote this as $\text{sign}(z_{h_{\ell}})$ for brevity. However, depending on the exact definition used there may be a difference for $z_{h_{\ell}} = 0$, as $\text{abs}'(0)$ is undefined. In practicality this doesn't matter much though, although theoretically it does mean that the expectation of this is theoretically undefined when $E[z_{h_{\ell}}] = 0$.

A.3 Gradient of NMU

In scalar notation the NMU is defined as:

$$z_{h_{\ell}} = \prod_{h_{\ell-1}=1}^{H_{\ell-1}} (W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}) \quad (25)$$

The gradient of the loss with respect to $W_{h_{\ell-1}, h_{\ell}}$ is fairly trivial. Note that every term but the one for $h_{\ell-1}$, is just a constant with respect to $W_{h_{\ell-1}, h_{\ell}}$. The product, expect the term for $h_{\ell-1}$ can be expressed as $\frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}}$. Using this fact, it becomes trivial to derive the gradient as:

$$\frac{\partial \mathcal{L}}{\partial w_{h_{\ell}, h_{\ell-1}}} = \frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \frac{\partial z_{h_{\ell}}}{\partial w_{h_{\ell}, h_{\ell-1}}} = \frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}} (z_{h_{\ell-1}} - 1) \quad (26)$$

Similarly, the gradient $\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}}$ which is essential in backpropagation can equally easily be derived as:

$$\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} = \sum_{h_{\ell}=1}^{H_{\ell}} \frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} = \sum_{h_{\ell}=1}^{H_{\ell}} \frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}} W_{h_{\ell-1}, h_{\ell}} \quad (27)$$

B Moments

B.1 Expectation and variance for weight matrix construction in NAC layers

The weight matrix construction in NAC, is defined in scalar notation as:

$$W_{h_{\ell}, h_{\ell-1}} = \tanh(\hat{W}_{h_{\ell}, h_{\ell-1}}) \sigma(\hat{M}_{h_{\ell}, h_{\ell-1}}) \quad (28)$$

Simplifying the notation of this, and re-expressing it using stochastic variables with uniform distributions this can be written as:

$$\begin{aligned}W &\sim \tanh(\hat{W}) \sigma(\hat{M}) \\ \hat{W} &\sim U[-r, r] \\ \hat{M} &\sim U[-r, r]\end{aligned}\quad (29)$$

Since $\tanh(\hat{W})$ is an odd-function and $E[\hat{W}] = 0$, deriving the expectation $E[W]$ is trivial.

$$E[W] = E[\tanh(\hat{W})] E[\sigma(\hat{M})] = 0 \cdot E[\sigma(\hat{M})] = 0 \quad (30)$$

The variance is more complicated, however as \hat{W} and \hat{M} are independent, it can be simplified to:

$$\text{Var}[W] = E[\tanh(\hat{W})^2] E[\sigma(\hat{M})^2] - E[\tanh(\hat{W})]^2 E[\sigma(\hat{M})]^2 = E[\tanh(\hat{W})^2] E[\sigma(\hat{M})^2] \quad (31)$$

These second moments can be analyzed independently. First for $E[\tanh(\hat{W})^2]$:

$$\begin{aligned}
E[\tanh(\hat{W})^2] &= \int_{-\infty}^{\infty} \tanh(x)^2 f_{U[-r,r]}(x) dx \\
&= \frac{1}{2r} \int_{-r}^r \tanh(x)^2 dx \\
&= \frac{1}{2r} \cdot 2 \cdot (r - \tanh(r)) \\
&= 1 - \frac{\tanh(r)}{r}
\end{aligned} \tag{32}$$

Then for $E[\tanh(\hat{M})^2]$:

$$\begin{aligned}
E[\sigma(\hat{M})^2] &= \int_{-\infty}^{\infty} \sigma(x)^2 f_{U[-r,r]}(x) dx \\
&= \frac{1}{2r} \int_{-r}^r \sigma(x)^2 dx \\
&= \frac{1}{2r} \left(r - \tanh\left(\frac{r}{2}\right) \right)
\end{aligned} \tag{33}$$

Finally this gives the variance:

$$\text{Var}[W] = \frac{1}{2r} \left(1 - \frac{\tanh(r)}{r} \right) \left(r - \tanh\left(\frac{r}{2}\right) \right) \tag{34}$$

B.2 Expectation and variance of NAC.

B.2.1 Forward pass

Assuming that each $z_{h_{\ell-1}}$ are independent the expectation can be simplified to:

$$\begin{aligned}
E[z_{h_{\ell}}] &= E \left[\exp \left(\sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_{\ell}, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon) \right) \right] \\
&= E \left[\prod_{h_{\ell-1}=1}^{H_{\ell-1}} \exp(W_{h_{\ell}, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon)) \right] \\
&= \prod_{h_{\ell-1}=1}^{H_{\ell-1}} E[\exp(W_{h_{\ell}, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon))] \\
&= E[\exp(W_{h_{\ell}, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon))]^{H_{\ell-1}} \\
&= E \left[(|z_{h_{\ell-1}}| + \epsilon)^{W_{h_{\ell}, h_{\ell-1}}} \right]^{H_{\ell-1}} \\
&= E \left[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}}) \right]^{H_{\ell-1}}
\end{aligned} \tag{35}$$

Here we define f as a non-linear transformation function of two independent stochastic variables:

$$f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}}) = (|z_{h_{\ell-1}}| + \epsilon)^{W_{h_{\ell}, h_{\ell-1}}} \tag{36}$$

We then take the second order taylor approximation of f , around $(E[z_{h_{\ell-1}}], E[W_{h_{\ell}, h_{\ell-1}}])$.

$$\begin{aligned}
E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})] &\approx E \left[\right. \\
&f(E[z_{h_{\ell-1}}], E[W_{h_{\ell}, h_{\ell-1}}]) \\
&+ \begin{bmatrix} z_{h_{\ell-1}} - E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} - E[W_{h_{\ell}, h_{\ell-1}}] \end{bmatrix}^T \begin{bmatrix} \frac{\partial f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial z_{h_{\ell-1}}} \\ \frac{\partial f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial W_{h_{\ell}, h_{\ell-1}}} \end{bmatrix} \left| \begin{array}{l} z_{h_{\ell-1}} = E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} = E[W_{h_{\ell}, h_{\ell-1}}] \end{array} \right. \right. \\
&+ \frac{1}{2} \begin{bmatrix} z_{h_{\ell-1}} - E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} - E[W_{h_{\ell}, h_{\ell-1}}] \end{bmatrix}^T \\
&\bullet \begin{bmatrix} \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial^2 z_{h_{\ell-1}}} & \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial z_{h_{\ell-1}} \partial W_{h_{\ell}, h_{\ell-1}}} \\ \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial z_{h_{\ell-1}} \partial W_{h_{\ell}, h_{\ell-1}}} & \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial^2 W_{h_{\ell}, h_{\ell-1}}} \end{bmatrix} \left| \begin{array}{l} z_{h_{\ell-1}} = E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} = E[W_{h_{\ell}, h_{\ell-1}}] \end{array} \right. \\
&\bullet \left. \begin{bmatrix} z_{h_{\ell-1}} - E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} - E[W_{h_{\ell}, h_{\ell-1}}] \end{bmatrix} \right] \quad (37)
\end{aligned}$$

Because $E[z_{h_{\ell-1}} - E[z_{h_{\ell-1}}]] = 0$, $E[W_{h_{\ell}, h_{\ell-1}} - E[W_{h_{\ell}, h_{\ell-1}}]] = 0$, and $Cov[z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}}] = 0$. This simplifies to:

$$\begin{aligned}
E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})] &\approx f(E[z_{h_{\ell-1}}], E[W_{h_{\ell}, h_{\ell-1}}]) \\
&+ \frac{1}{2} Var \begin{bmatrix} z_{h_{\ell-1}} \\ W_{h_{\ell}, h_{\ell-1}} \end{bmatrix}^T \begin{bmatrix} \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial^2 z_{h_{\ell-1}}} \\ \frac{\partial^2 f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})}{\partial^2 W_{h_{\ell}, h_{\ell-1}}} \end{bmatrix} \left| \begin{array}{l} z_{h_{\ell-1}} = E[z_{h_{\ell-1}}] \\ W_{h_{\ell}, h_{\ell-1}} = E[W_{h_{\ell}, h_{\ell-1}}] \end{array} \right. \quad (38)
\end{aligned}$$

Inserting the derivatives and computing the inner products yields:

$$\begin{aligned}
E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})] &\approx (|E[z_{h_{\ell-1}}]| + \epsilon)^{E[W_{h_{\ell}, h_{\ell-1}}]} \\
&+ \frac{1}{2} Var[z_{h_{\ell-1}}] (|E[z_{h_{\ell-1}}]| + \epsilon)^{E[W_{h_{\ell}, h_{\ell-1}}]-2} E[W_{h_{\ell}, h_{\ell-1}}] (E[W_{h_{\ell}, h_{\ell-1}}] - 1) \\
&+ \frac{1}{2} Var[W_{h_{\ell}, h_{\ell-1}}] (|E[z_{h_{\ell-1}}]| + \epsilon)^{E[W_{h_{\ell}, h_{\ell-1}}]} \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \\
&= 1 + \frac{1}{2} Var[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \quad (39)
\end{aligned}$$

This gives the final expectation:

$$\begin{aligned}
E[z_{h_{\ell}}] &= E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})]^{H_{\ell-1}} \\
&\approx \left(1 + \frac{1}{2} Var[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \right)^{H_{\ell-1}} \quad (40)
\end{aligned}$$

As this expectation is of particular interest, we evaluate the error of the approximation, where $W_{h_{\ell}, h_{\ell-1}} \sim U[-r_w, r_w]$ and $z_{h_{\ell-1}} \sim U[0, r_z]$. These distributions are what is used in the simple function task is done. The result is seen in figure 9.

The variance can be derived using the same assumptions about expectation and no correlation.

$$\begin{aligned}
Var[z_{h_{\ell}}] &= E[z_{h_{\ell}}^2] - E[z_{h_{\ell}}]^2 \\
&= E \left[\prod_{h_{\ell-1}=1}^{H_{\ell-1}} (|z_{h_{\ell-1}}| + \epsilon)^{2 \cdot W_{h_{\ell}, h_{\ell-1}}} \right] - E \left[\prod_{h_{\ell-1}=1}^{H_{\ell-1}} (|z_{h_{\ell-1}}| + \epsilon)^{W_{h_{\ell}, h_{\ell-1}}} \right]^2 \\
&= E[f(z_{h_{\ell-1}}, 2 \cdot W_{h_{\ell}, h_{\ell-1}})]^{H_{\ell-1}} - E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})]^{2 \cdot H_{\ell-1}} \quad (41)
\end{aligned}$$

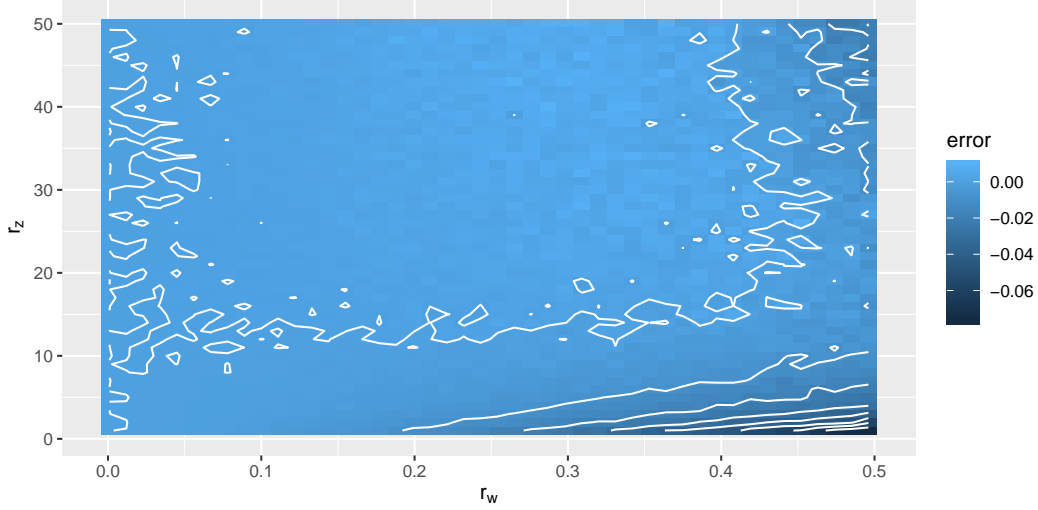


Figure 9: Error between theoretical approximation and the numerical approximation estimated by random sampling of 100000 observations at each combination of r_z and r_w .

We already have from the expectation result that:

$$E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})] \approx 1 + \frac{1}{2} \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \quad (42)$$

By substitution of variable we have that:

$$\begin{aligned} E[f(z_{h_{\ell-1}}, 2 \cdot W_{h_{\ell}, h_{\ell-1}})] &\approx 1 + \frac{1}{2} \text{Var}[2 \cdot W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \\ &\approx 1 + 2 \cdot \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \end{aligned} \quad (43)$$

This gives the variance:

$$\begin{aligned} \text{Var}[z_{h_{\ell}}] &= E[f(z_{h_{\ell-1}}, 2 \cdot W_{h_{\ell}, h_{\ell-1}})]^{H_{\ell-1}} - E[f(z_{h_{\ell-1}}, W_{h_{\ell}, h_{\ell-1}})]^{2 \cdot H_{\ell-1}} \\ &\approx (1 + 2 \cdot \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2)^{H_{\ell-1}} \\ &\quad - \left(1 + \frac{1}{2} \cdot \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2\right)^{2 \cdot H_{\ell-1}} \end{aligned} \quad (44)$$

B.2.2 Backward pass

The expectation of the backpropagation term:

$$E[\delta_{h_{\ell}}] = E\left[\sum_{h_{\ell+1}=1}^{H_{\ell+1}} \delta_{h_{\ell+1}} \frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}}\right] = H_{\ell+1} E[\delta_{h_{\ell+1}}] E\left[\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}}\right] \quad (45)$$

Where we have that:

$$E\left[\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}}\right] = E[h_{\ell+1}] E[W_{h_{\ell+1}, h_{\ell}}] E\left[\frac{\text{abs}'(z_{h_{\ell}})}{|z| + \epsilon}\right] = E[m_{h_{\ell+1}}] \cdot 0 \cdot E\left[\frac{\text{abs}'(z_{h_{\ell}})}{|z| + \epsilon}\right] = 0 \quad (46)$$

Deriving the variance is more complicated as:

$$\text{Var}\left[\frac{\partial m_{h_{\ell+1}}}{\partial z_{h_{\ell}}}\right] = \text{Var}\left[m_{h_{\ell+1}} W_{h_{\ell+1}, h_{\ell}} \frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon}\right] \quad (47)$$

Assuming independence between each term this can be simplified to as:

$$\begin{aligned}
\text{Var} \left[\frac{\partial z_{h_{\ell+1}}}{\partial z_{h_{\ell}}} \right] &= E[z_{h_{\ell+1}}^2] E[W_{h_{\ell+1}, h_{\ell}}^2] E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \right)^2 \right] \\
&\quad - E[z_{h_{\ell+1}}]^2 E[W_{h_{\ell+1}, h_{\ell}}]^2 E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \right)^2 \right] \\
&= E[z_{h_{\ell+1}}^2] \text{Var}[W_{h_{\ell+1}, h_{\ell}}] E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \right)^2 \right] \\
&\quad - E[z_{h_{\ell+1}}]^2 \cdot 0 \cdot E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \right)^2 \right] \\
&= E[z_{h_{\ell+1}}^2] \text{Var}[W_{h_{\ell+1}, h_{\ell}}] E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z_{h_{\ell}}| + \epsilon} \right)^2 \right]
\end{aligned} \tag{48}$$

Using Taylor approximation around $E[z_{h_{\ell}}]$ we have:

$$\begin{aligned}
E \left[\left(\frac{\text{abs}'(z_{h_{\ell}})}{|z| + \epsilon} \right)^2 \right] &\approx \frac{1}{(|E[z_{h_{\ell}}]| + \epsilon)^2} + \frac{1}{2} \frac{6}{(|E[z_{h_{\ell}}]| + \epsilon)^4} \text{Var}[z_{h_{\ell}}] \\
&= \frac{1}{(|E[z_{h_{\ell}}]| + \epsilon)^2} + \frac{3}{(|E[z_{h_{\ell}}]| + \epsilon)^4} \text{Var}[z_{h_{\ell}}]
\end{aligned} \tag{49}$$

Also reusing the result for $E[z_{h_{\ell}}^2]$ from earlier the variance can be expressed as:

$$\begin{aligned}
\text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} \right] &\approx \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} (1 + 2 \cdot \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2)^{H_{\ell-1}} \\
&\quad \cdot \text{Var}[W_{h_{\ell}, h_{\ell-1}}] \left(\frac{1}{(|E[z_{h_{\ell-1}}]| + \epsilon)^2} + \frac{3}{(|E[z_{h_{\ell-1}}]| + \epsilon)^4} \text{Var}[z_{h_{\ell-1}}] \right)
\end{aligned} \tag{50}$$

B.3 Expectation and variance of NMu

B.3.1 Forward pass

Assuming that each $z_{h_{\ell-1}}$ are independent, that $E[z_{h_{\ell-1}}] = 0$, and that $E[W_{h_{\ell-1}, h_{\ell}}] = 1/2$ the expectation is:

$$\begin{aligned}
E[z_{h_{\ell}}] &\approx E \left[\prod_{h_{\ell-1}=1}^{H_{\ell-1}} (W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}) \right] \\
&\approx E [W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}]^{H_{\ell-1}} \\
&\approx (E[W_{h_{\ell-1}, h_{\ell}}] E[z_{h_{\ell-1}}] + 1 - E[W_{h_{\ell-1}, h_{\ell}}])^{H_{\ell-1}} \\
&\approx \left(\frac{1}{2} \cdot 0 + 1 - \frac{1}{2} \right)^{H_{\ell-1}} \\
&\approx \left(\frac{1}{2} \right)^{H_{\ell-1}}
\end{aligned} \tag{51}$$

Using the same assumptions for the variance one gets:

$$\begin{aligned}
\text{Var}[z_{h_\ell}] &= E[z_{h_\ell}^2] - E[z_{h_\ell}]^2 \\
&\approx E[z_{h_\ell}^2] - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&\approx E \left[\prod_{h_{\ell-1}=1}^{H_{\ell-1}} (W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell})^2 \right] - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&\approx E[(W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell})^2]^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&\approx \left(E[W_{h_{\ell-1}, h_\ell}^2] E[z_{h_{\ell-1}}^2] - 2E[W_{h_{\ell-1}, h_\ell}] E[z_{h_{\ell-1}}] \right. \\
&\quad \left. + E[W_{h_{\ell-1}, h_\ell}^2] + 2E[W_{h_{\ell-1}, h_\ell}] E[z_{h_{\ell-1}}] \right. \\
&\quad \left. - 2E[W_{h_{\ell-1}, h_\ell}] + 1 \right)^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \tag{52} \\
&\approx \left(E[W_{h_{\ell-1}, h_\ell}^2] E[z_{h_{\ell-1}}^2] + E[W_{h_{\ell-1}, h_\ell}^2] \right. \\
&\quad \left. - 2E[W_{h_{\ell-1}, h_\ell}] + 1 \right)^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&= \left(E[W_{h_{\ell-1}, h_\ell}^2] (E[z_{h_{\ell-1}}^2] + 1) \right)^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&\approx ((\text{Var}[W_{h_{\ell-1}, h_\ell}] + E[W_{h_{\ell-1}, h_\ell}]^2) (\text{Var}[z_{h_{\ell-1}}] + 1))^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}} \\
&= \left(\text{Var}[W_{h_{\ell-1}, h_\ell}] + \frac{1}{4} \right)^{H_{\ell-1}} (\text{Var}[z_{h_{\ell-1}}] + 1)^{H_{\ell-1}} - \left(\frac{1}{2}\right)^{2 \cdot H_{\ell-1}}
\end{aligned}$$

B.3.2 Backward pass

For the backward pass the expectation can using the same assumptions be derived to:

$$\begin{aligned}
E \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} \right] &= H_\ell E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \frac{\partial z_{h_\ell}}{\partial z_{h_{\ell-1}}} \right] \\
&= H_\ell E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \right] E \left[\frac{\partial z_{h_\ell}}{\partial z_{h_{\ell-1}}} \right] \\
&= H_\ell E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \right] E \left[\frac{z_{h_\ell}}{W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell}} W_{h_{\ell-1}, h_\ell} \right] \\
&= H_\ell E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \right] E \left[\frac{z_{h_\ell}}{W_{h_{\ell-1}, h_\ell} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_\ell}} \right] E[W_{h_{\ell-1}, h_\ell}] \tag{53} \\
&= H_\ell E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \right] \left(\frac{1}{2}\right)^{H_{\ell-1}-1} \frac{1}{2} \\
&= E \left[\frac{\partial \mathcal{L}}{\partial z_{h_\ell}} \right] H_\ell \left(\frac{1}{2}\right)^{H_{\ell-1}} \\
&\approx 0 \cdot H_\ell \cdot \left(\frac{1}{2}\right)^{H_{\ell-1}} \\
&= 0
\end{aligned}$$

And finally the variance for the backward pass is derived using the same assumptions.

$$\begin{aligned}
\text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell-1}}} \right] &= H_{\ell} \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} \right] \\
&= H_{\ell} \left(\text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] E \left[\frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} \right]^2 + E \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right]^2 \text{Var} \left[\frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} \right] \right. \\
&\quad \left. + \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] \text{Var} \left[\frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} \right] \right) \\
&\approx \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} \text{Var} \left[\frac{\partial z_{h_{\ell}}}{\partial z_{h_{\ell-1}}} \right] \\
&\approx \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} \left(E \left[\left(\frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}} \right)^2 \right] E[W_{h_{\ell-1}, h_{\ell}}^2] \right. \\
&\quad \left. - E \left[\frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}} \right]^2 E[W_{h_{\ell-1}, h_{\ell}}]^2 \right) \\
&\approx \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} \left(E \left[\left(\frac{z_{h_{\ell}}}{W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}}} \right)^2 \right] E[W_{h_{\ell-1}, h_{\ell}}^2] \right. \\
&\quad \left. - \left(\frac{1}{2} \right)^{2(H_{\ell-1}-1)} \left(\frac{1}{2} \right)^2 \right) \\
&\approx \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} \left(\left(\left(\text{Var}[W_{h_{\ell-1}, h_{\ell}}] + \frac{1}{4} \right) (\text{Var}[z_{h_{\ell-1}}] + 1) \right)^{H_{\ell-1}-1} \right. \\
&\quad \left. \cdot \left(\text{Var}[W_{h_{\ell-1}, h_{\ell}}] + \frac{1}{4} \right) - \left(\frac{1}{2} \right)^{2 \cdot H_{\ell-1}} \right) \\
&= \text{Var} \left[\frac{\partial \mathcal{L}}{\partial z_{h_{\ell}}} \right] H_{\ell} \left(\left(\text{Var}[W_{h_{\ell-1}, h_{\ell}}] + \frac{1}{4} \right)^{H_{\ell-1}} (\text{Var}[z_{h_{\ell-1}}] + 1)^{H_{\ell-1}-1} \right. \\
&\quad \left. - \left(\frac{1}{2} \right)^{2 \cdot H_{\ell-1}} \right)
\end{aligned} \tag{54}$$

B.3.3 Initialization

The expectation of $W_{h_{\ell-1}, h_{\ell}}$ should be $E[W_{h_{\ell-1}, h_{\ell}}] = \frac{1}{2}$. Using the variance approximations found, the variance should be according to the forward pass:

$$\text{Var}[W_{h_{\ell-1}, h_{\ell}}] = ((1 + \text{Var}[z_{h_{\ell}}])^{-H_{\ell-1}} \text{Var}[z_{h_{\ell}}] + (4 + 4\text{Var}[z_{h_{\ell}}])^{-H_{\ell-1}})^{\frac{1}{H_{\ell-1}}} - \frac{1}{4} \tag{55}$$

And according to the backward pass it should be:

$$\text{Var}[W_{h_{\ell-1}, h_{\ell}}] = (H_{\ell}(1 + \text{Var}[z_{h_{\ell-1}}])(4 + 4\text{Var}[z_{h_{\ell-1}}])^{-H_{\ell-1}} + (1 + \text{Var}[z_{h_{\ell-1}}])^{1-H_{\ell-1}})^{\frac{1}{H_{\ell-1}}} - \frac{1}{4} \tag{56}$$

These are both dependent on the input variance. If the input variance is known then optimal initialization is possible. However, as this is often not the case one can perhaps assume that $\text{Var}[z_{h_{\ell-1}}] = 1$. This is not an unreasonable assumption in many cases, as there may either be a normalization layer somewhere or the input is normalized. If unit variance is assumed, one gets from the forward pass:

$$\text{Var}[W_{h_{\ell-1}, h_{\ell}}] = (2^{-H_{\ell-1}} + 8^{-H_{\ell-1}})^{\frac{1}{H_{\ell-1}}} - \frac{1}{4} = \frac{1}{8} \left((4^{H_{\ell-1}} + 1)^{H_{\ell-1}} - 2 \right) \tag{57}$$

And from the backward pass:

$$\text{Var}[W_{h_{\ell-1}, h_{\ell}}] = (2H_{\ell}8^{-H_{\ell-1}} + 2^{1-H_{\ell-1}})^{\frac{1}{H}} - \frac{1}{4} \quad (58)$$

The variance requirement for the backward pass is hard to satisfy, as this is dependent on two variables. However, the variance requirement from the forward pass quickly $\text{Var}[W_{h_{\ell-1}, h_{\ell}}] = \frac{1}{4}$ may be a reasonable initialization.

C Simple function task

C.1 Dataset generation

All datasets in the simple function task experiments are generated using the following algorithm:

Algorithm 1 Dataset sampling algorithm

```

1: function DATASET(OP( $\cdot, \cdot$ ) : Operation,  $i$  : InputSize,  $s$  : SubsetRatio,  $o$  : OverlapRatio,
    $R$  : Range)
2:    $\mathbf{x} \leftarrow \text{UNIFORM}(R_{\text{lower}}, R_{\text{upper}}, i)$                                  $\triangleright$  Sample  $i$  elements uniformly
3:    $k \leftarrow \text{UNIFORM}(0, 1 - 2s - o)$                                         $\triangleright$  Sample offset
4:    $a \leftarrow \text{SUM}(\mathbf{x}[ik : i(k + s)])$                                         $\triangleright$  Create sum  $a$  from subset
5:    $b \leftarrow \text{SUM}(\mathbf{x}[i(k + s - o) : i(k + 2s - 0)])$                     $\triangleright$  Create sum  $b$  from subset
6:    $t \leftarrow \text{OP}(a, b)$                                                      $\triangleright$  Perform operation on  $a$  and  $b$ 
7:   return  $x, t$ 

```
