

On the vanishing and exploding gradient problem in Gated Recurrent Units

Alexander Rehmer, Andreas Kroll

Department of Measurement and Control, Institute for System Analytics and Control, Faculty of Mechanical Engineering, University of Kassel, Germany, (e-mail: {alexander.rehmer, andreas.kroll}@mrt.uni-kassel.de)

Abstract: Recurrent Neural Networks are applied in areas such as speech recognition, natural language and video processing, and the identification of nonlinear state space models. Conventional Recurrent Neural Networks, e.g. the Elman Network, are hard to train. A more recently developed class of recurrent neural networks, so-called Gated Units, outperform their counterparts on virtually every task. This paper aims to provide additional insights into the differences between RNNs and Gated Units in order to explain the superior performance of gated recurrent units. It is argued, that Gated Units are easier to optimize not because they solve the vanishing gradient problem, but because they circumvent the emergence of large local gradients.

Copyright © 2020 The Authors. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0>)

Keywords: Nonlinear system identification, Recurrent Neural Networks, Gated Recurrent Units.

1. INTRODUCTION

Gated Units, such as the Long Short-Term Memory (LSTM) and the Gated Recurrent Unit (GRU) were originally developed to overcome the *vanishing gradient problem*, which occurs in the Elman Recurrent Neural Network (RNN) (Pascanu et al., 2012). They have since outperformed RNNs on a number of tasks, such as natural language, speech and video processing (Jordan et al., 2019) and recently also on a nonlinear system identification task (Rehmer and Kroll, 2019). However, it will be shown, that the gradient also vanishes in Gated Units, other still unaccounted for mechanisms have to be responsible for their success. Pascanu et al. (2012) show, that small changes in the parameters θ of the RNN can lead to drastic changes in the dynamic behavior of the system, when crossing certain critical bifurcation points. This in turn results in a huge change in the evolution of the hidden state \hat{x}_k , which leads to a locally large, or *exploding*, gradient of the loss function. In this paper the GRU will be examined and compared to the RNN with the purpose to provide an alternative explanation to why GRUs outperform RNNs. First, it will be shown, that the gradient of the GRU is in fact smaller than that of the RNN, at least for the parameterizations considered in this paper, although GRUs were originally designed to solve the vanishing gradient problem. Secondly, it will be shown, that GRUs are not only capable to represent highly nonlinear dynamics, but are also able to represent approximately linear dynamics via a number of different parameterizations. Since a linear model is always a good first guess, the easy accessibility of different parameterizations that produce linear dynam-

ics makes the GRU less sensitive to its initial choice of parameters and thus simplifies the optimization problem. In the end GRU and RNNs will be compared on a simple academic example and on a real nonlinear identification task.

2. RECURRENT NEURAL NETWORKS

In this section the Simple Recurrent Neural Network (RNN), also known as Elman Network, and the Gated Recurrent Unit (GRU) will be introduced.

2.1 Simple Recurrent Neural Network

The RNN as depicted in figure 1 is a straightforward realization of a nonlinear state space model (Nelles, 2001). It consists of one hidden recurrent layer with nonlinear activation function f_h , which aims to approximate the state equation, as well as one hidden feedforward layer f_g and one linear output layer, which together aim to approximate the output equation. For simplicity of notation, the linear output layer is omitted in the following equations and figures:

$$\begin{aligned}\hat{x}_{k+1} &= f_h(W_x \hat{x}_k + W_u u_k + b_h), \\ \hat{y}_k &= f_g(W_y \hat{x}_k + b_g),\end{aligned}\quad (1)$$

with $\hat{x}_k \in \mathbb{R}^{n \times 1}$, $\hat{y}_k \in \mathbb{R}^{m \times 1}$, $u_k \in \mathbb{R}^{l \times 1}$, $W_x \in \mathbb{R}^{n \times n}$, $W_u \in \mathbb{R}^{n \times l}$, $b_h \in \mathbb{R}^{n \times 1}$, $W_y \in \mathbb{R}^{m \times n}$, $b_g \in \mathbb{R}^{m \times 1}$ and $f_h: \mathbb{R}^{n \times 1} \rightarrow \mathbb{R}^{n \times 1}$, $f_g: \mathbb{R}^{m \times 1} \rightarrow \mathbb{R}^{m \times 1}$. Usually $\tanh(\cdot)$ is employed as nonlinear activation function. When training an RNN, the recurrent model is unfolded over the whole training sequence of length N , and the gradient of the loss function \mathcal{L} with respect to the model parameters θ is calculated. As a consequence of the feedback, the gradient of the error

$$e_k = (\hat{y}_k - y_k^{\text{target}}) \quad (2)$$

* Sponsor and financial support acknowledgment goes here. Paper titles should be written in uppercase and lowercase letters, not all uppercase.

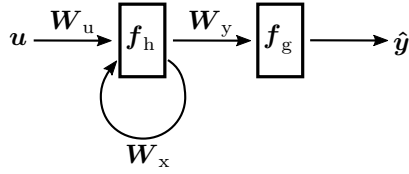


Fig. 1. Representation of the Elman Network: Layers of neurons are represented as rectangles, connections between layers represent fully connected layers.

at time step k with respect to the model parameters $\theta = \{W_x, W_u, W_y, b_h, b_g\}$ depends on the previous state \hat{x}_{k-1} , which depends again on the model parameters:

$$\frac{\partial e_k}{\partial \theta} = \frac{\partial e_k}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \hat{x}_k} \left[\frac{\partial \hat{x}_k}{\partial \theta} + \frac{\partial \hat{x}_k}{\partial \hat{x}_{k-1}} \frac{\partial \hat{x}_{k-1}}{\partial \theta} \right] \quad (3)$$

For example, the gradient of the hidden state \hat{x}_k with respect to W_x is

$$\frac{\partial \hat{x}_k}{\partial W_x} = \sum_{\tau=1}^N \hat{x}_{k-\tau} f_h'^{(k-\tau+1)}(\cdot) \prod_{\beta} f_h'^{(k-\beta)}(\cdot) W_x^{(k-\beta)} \quad (4)$$

$$\beta = \tau - 2, \tau - 3, \dots \forall \beta \geq 0.$$

High indices in brackets indicate the particular time step. The product term in (4), which also appears when computing the gradient with respect to the other parameters, decreases exponentially with τ , if $|f_h' \cdot \rho(W)| < 1$, where $\rho(W_x)$ is the spectral radius of W_x . Essentially, backpropagating an error one time step involves a multiplication of the state with a derivative that is possibly smaller than one and a matrix whose spectral radius is possibly smaller than one. Hence, the gradient vanishes after a certain amount of time steps. In the Machine Learning community it is argued that the vanishing gradient prevents learning of so-called *long-term dependencies* in acceptable time (Hochreiter and Schmidhuber, 1997; Goodfellow et al., 2016), i.e. when huge time lags exist between input u_k and output \hat{y}_k . Gated recurrent units, like LSTM and GRU were developed to solve this problem and have since then outperformed classical RNNs on virtually any task. However, it can be shown that the gradient also vanishes in gated recurrent units. Additionally, the vanishing of the gradient over time is a desirable property. In most systems, the influence of a previous state $\hat{x}_{k-\tau}$ on the current state \hat{x}_k decreases over time. Unless one wants to design a marginally stable or unstable system, e.g. when performing tasks like unbounded counting, or when dealing with large dead times, the vanishing gradient has no negative effect on the optimization procedure.

2.2 The Gated Recurrent Unit

Gated Recurrent Unit (GRU) (Cho et al., 2014) is besides LSTM the most often applied architecture of gated recurrent units. The general concept of gated recurrent units is to manipulate the state \hat{x}_k through the addition or multiplication of the activations of so called gates, see figure 2. Gates are almost exclusively one-layered neural networks with nonlinear sigmoid activation functions. The state equation of the GRU is

$$\hat{x}_{k+1} = f_z \odot \hat{x}_k + (1 - f_z) \odot f_c(\tilde{x}_k) \quad (5)$$

with $\tilde{x}_k = f_r \odot \hat{x}_k$. The operator \odot denotes the Hadamard product. The activations of the so-called gate reset gate f_r , update gate f_z and the output gate f_c are given by

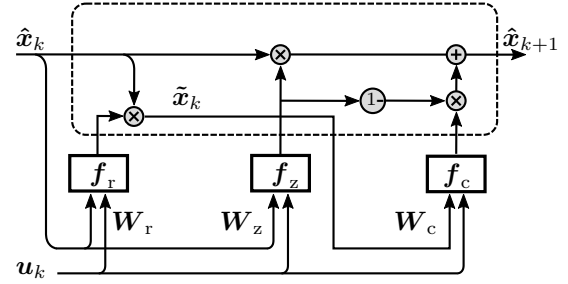


Fig. 2. The Gated Recurrent Unit (GRU). Gates are depicted as rectangles with their respective activation functions.

$$\begin{aligned} f_r &= \sigma(W_r \cdot [\hat{x}_k, u_k] + b_r), \\ f_z &= \sigma(W_z \cdot [\hat{x}_k, u_k] + b_z), \\ f_c &= \tanh(W_c \cdot [\hat{x}_k, u_k] + b_c), \end{aligned} \quad (6)$$

where $W_r, W_z, W_c \in \mathbb{R}^{n \times n+l}$, $b_r, b_z, b_c \in \mathbb{R}^{n \times 1}$ and $f_r, f_z, f_c : \mathbb{R}^{n \times 1} \rightarrow \mathbb{R}^{n \times 1}$. $\sigma(\cdot)$ denotes the logistic function. In order to map the states estimated by the GRU to the output, the GRU has to be equipped either with an output layer, as the RNN, or with an output gate, as the LSTM.

3. GRADIENT OF THE STATE EQUATIONS

In this section, the gradients of the state equations of RNN (1) and GRU (5) w.r.t. their parameters will be compared to each other. In the cases examined the gradient of the GRU is, somewhat surprisingly, at most as large as that of the RNN, but usually smaller.

In order to allow for an easily interpretable visualization, the analysis will be restricted to one dimensional and autonomous systems. Also, the GRU will be simplified by eliminating the reset gate f_r from (6), such that $\tilde{x}_k = \hat{x}_k$. Taking the gradient of the RNN's state equation in (1) w.r.t. w_x yields

$$\frac{\partial \hat{x}_{k+1}}{\partial w_x} = (\hat{x}_k + w_x \frac{\partial \hat{x}_k}{\partial w_x}) \cdot \tanh'(w_x \hat{x}_k + b_x). \quad (7)$$

The gradient of the GRUs state equation (5) with respect to w_z is

$$\begin{aligned} \frac{\partial \hat{x}_{k+1}}{\partial w_z} &= (1 - \tanh(\hat{x}_k; \theta_z)) \sigma'(\hat{x}_k; \theta_z) \left(\hat{x}_k + w_z \frac{\partial \hat{x}_k}{\partial w_z} \right) \\ &\quad + (1 - \sigma(\hat{x}_k; \theta_z)) \tanh'(\hat{x}_k; \theta_c) \frac{\partial \hat{x}_k}{\partial w_z}, \end{aligned} \quad (8)$$

and with respect to w_c

$$\begin{aligned} \frac{\partial \hat{x}_{k+1}}{\partial w_c} &= \sigma'(\hat{x}_k; \theta_z) \frac{\partial \hat{x}_k}{\partial w_c} (\hat{x}_k - \tanh(\hat{x}_k; \theta_c)) \\ &\quad + \sigma(\hat{x}_k; \theta_z) \frac{\partial \hat{x}_k}{\partial w_c} \\ &\quad + (1 - \sigma(\hat{x}_k; \theta_z)) \tanh'(\hat{x}_k; \theta_c) \left(\hat{x}_k + w_c \frac{\partial \hat{x}_k}{\partial w_c} \right). \end{aligned} \quad (9)$$

For convenience of notation θ_z and θ_c denote the parameters of f_z and f_c respectively, i.e. $\theta_z = [w_z, b_z]$ and $\theta_c = [w_c, b_c]$.

It is cumbersome to write down the gradients in (8) and (9) for an arbitrary number of time steps, as done for the

RNN in (4). However, it already becomes apparent from (8) and (9), that each step of backpropagation involves the multiplication of the previous state with the output of a sigmoid function, either $\sigma(\cdot)$ or $\tanh(\cdot)$, and with one of its derivatives. The outputs of the sigmoid functions is always smaller than one. The derivative of the $\tanh(\cdot)$ is smaller than one, if $|\mathbf{W}_c|_1 < 1$ and the derivative of the logistic function is smaller than one, if $|\mathbf{W}_z|_1 < 4$. This means, the gradient of the GRU is at least as likely to vanish as the gradient of the RNN.

In order to examine the gradients further it is assumed, that $|w_x|, |w_z|, |w_c| \leq 1$. The derivative of the logistic function $\sigma(\cdot)$ is therefore bounded in the interval $[0, \frac{1}{4}]$ and the derivative of the $\tanh(\cdot)$ bounded in the interval $[0, 1]$. Under these assumptions the gradient of the RNN is at most

$$\frac{\partial \hat{x}_{k+1}}{\partial w_x} \leq (\hat{x}_k + w_x \frac{\partial \hat{x}_k}{\partial w_x}). \quad (10)$$

For the gradient of the GRU two extreme cases will be examined and compared to (10).

If the update gate is fully closed, i.e. $\sigma(\hat{x}_k; \theta_z) = 0$, (8) and (9) become

$$\begin{aligned} \frac{\partial \hat{x}_{k+1}}{\partial w_z} &= \tanh'(\hat{x}_k; \theta_c) \frac{\partial \hat{x}_k}{\partial w_z} \leq \frac{\partial \hat{x}_k}{\partial w_z} \\ \frac{\partial \hat{x}_{k+1}}{\partial w_c} &= \tanh'(\hat{x}_k; \theta_c) (\hat{x}_k + w_c \frac{\partial \hat{x}_k}{\partial w_c}) \\ &\leq \hat{x}_k + w_c \frac{\partial \hat{x}_k}{\partial w_c} \end{aligned} \quad (11)$$

It follows, that if the update gate is fully closed, the GRU becomes an ordinary RNN with the same properties regarding its gradient.

If the update gate is fully open on the other hand, i.e. $\sigma(\hat{x}_k; \theta_z) = 1$, (8) and (9) become

$$\begin{aligned} \frac{\partial \hat{x}_{k+1}}{\partial w_z} &= 0, \\ \frac{\partial \hat{x}_{k+1}}{\partial w_c} &= \frac{\partial \hat{x}_k}{\partial w_c}. \end{aligned} \quad (12)$$

Essentially, the previous state \hat{x}_k is just passed over to the new state \hat{x}_{k+1} without any modification, so there is no gradient w.r.t. the parameters but that from the previous time-step.

One can examine the behavior of (8) and (9) in between these extreme cases by replacing $\tanh'(\cdot) = 1 - \tanh^2(\cdot)$ and $\sigma'(\cdot) = \sigma(\cdot)(1 - \sigma(\cdot))$. (8) and (9) then each are a cubic polynomial in $\sigma(\hat{x}_k, \theta_z)$ and $\tanh(\hat{x}_k, \theta_c)$. Their graphs are plotted in figure 3 for $w_z = w_c = 1$, which produces the largest gradient possible under the assumptions made. Also we set $\frac{\partial \hat{x}_k}{\partial w_z} = \frac{\partial \hat{x}_k}{\partial w_c} = 1$ and $\hat{x}_k = 1$.

Figure 3 confirms, that the largest gradients occurs, when the GRU's update gate is fully closed, i.e. $\sigma(\hat{x}_k, \theta_z) = 0$, and the GRU hence becomes an RNN. For all other configurations, the gradient is in fact smaller.

In contrast to the point made by the *vanishing gradient* argument, that if the gradient vanishes after some time-steps, optimization might take prohibitively long. We argue that a smaller gradient of the state equation w.r.t. its parameters is beneficial when training recurrent neural networks. If the gradient is large, as is the case with the RNN, a small change in parameters will lead to a huge change in the state trajectory. Since every time step, the previous state is again fed to the recurrent network, these

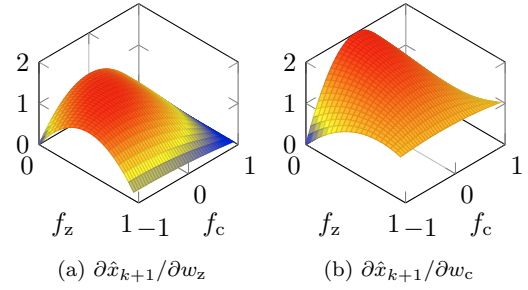


Fig. 3. Gradients of the GRU w.r.t. w_z and w_c depending on the activation of $f_z = \sigma(\hat{x}_k, \theta_z)$ and $f_c = \tanh(\hat{x}_k, \theta_c)$ for $w_z = w_c = \frac{\partial \hat{x}_k}{\partial w_z} = \frac{\partial \hat{x}_k}{\partial w_c} = \hat{x}_k = 1$

changes accumulate quickly, resulting in vastly different evolutions of the state. The effect on the loss function are locally huge gradients, which make gradient-based optimization very difficult. Compared to the RNN, the GRUs gradient of the state equation is almost always smaller, which means similar parameters produce similar trajectories of the state, leading to a smoother loss function without huge local gradients. This will be illustrated via simple examples in section 5.

4. EFFECT OF DIFFERENT PARAMETERIZATIONS ON THE STATE EQUATION

In addition to the smaller gradient of the GRU, the rather special structure of its state equation also benefits gradient based optimization, especially when identifying technical systems.

In this section, the same assumptions and restrictions apply as in the previous section, i.e. only one dimensional and autonomous systems are investigated and the GRUs reset gate f_r is neglected.

Without reset gate, (5) is the sum of the identity function, weighted with a logistic function $\sigma(\cdot)$, and a $\tanh(\cdot)$, weighted by the residual $1 - \sigma(\cdot)$. Fig. 4 visually decomposes (5) into its constituent parts.

This superposition of multiple functions enables the GRU to approximate quite complex nonlinearities with comparatively few parameters.

For example, if $b_z \rightarrow -\infty$, (5) converges to a tanh parameterized by w_c and b_c . For $w_z, b_c \rightarrow 0$, $w_c \rightarrow +\infty$ and $b_z \rightarrow -\infty$, (5) converges to the binary step function, while for small negative b_z it represents a kind of leaky binary step, where the slope of the tails is \tilde{b}_z , compare Figure 5:

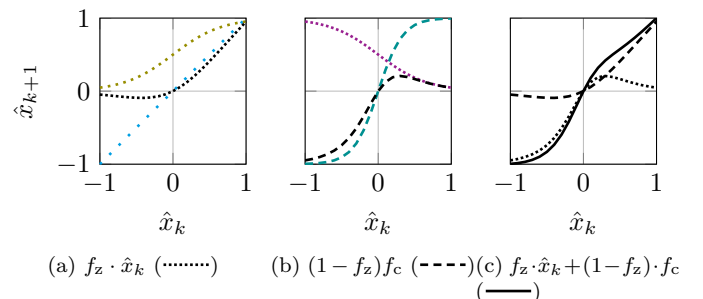


Fig. 4. Graphical decomposition of the GRU's state equation. \hat{x}_k (\cdots), f_z (\cdots), f_c ($---$), $(1 - f_z)$ (\cdots)

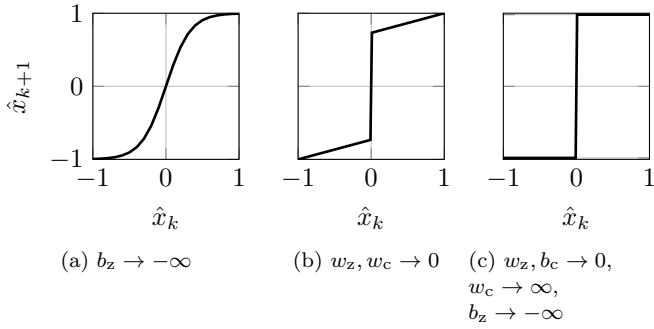


Fig. 5. Parameterizations for which GRU converges towards a tanh, leaky binary step or binary step activation function.

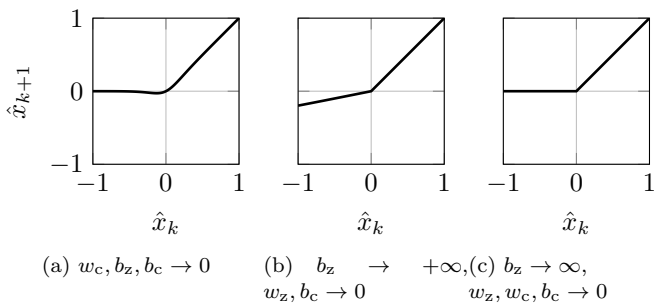


Fig. 6. Parameterizations for which GRU converges to variants of the ReLU activation function.

$$\tilde{b}_z = \sigma(b_z)$$

(5) also converges to the ReLU activation function for $w_z \rightarrow \pm\infty$ and $w_c, b_z, b_c \rightarrow \pm 0$ and is also able to represent a leaky ReLU with tail slope of \tilde{w}_c :

$$\tilde{w}_c = \tanh(w_c)$$

for small w_c and the swish activation for small w_z . These configurations are depicted in Figure 6.

The fact that the GRUs state equation converges to some of these functions when certain parameters become larger has an important effect on the loss function: Assume one tries to learn a state space equation in form of a tanh. Although the RNN would be the right choice, because its in the right model class, the optimization problem is rather difficult. The only solution lies in a steep valley, which means small deviations from that solution produce large losses. It would take infinitely small steps (as soon as you are in the valley, which one cannot know) to reach that solution. The GRU on the other hand offers an infinite amount of solutions for large b_z , which corresponds to an infinite plane with small slope in the parameter space. One only has to move along that plane to approximate the system increasingly better, which is easier than diving into a steep valley.

Additionally, the GRU is able to approximate the identity function and linear functions in general via a number of different parameterizations.: For $b_z \rightarrow +\infty$ (5) converges to the identity function, since f_z is approaching one. For $w_z, w_c \rightarrow 0$ the result is also a linear function, but its slope is determined by \tilde{b}_z and its axis intercept by \tilde{b}_c , see Figure 6:

$$\tilde{b}_c = \tanh(b_c)$$

The ability to converge to a linear model when approach-

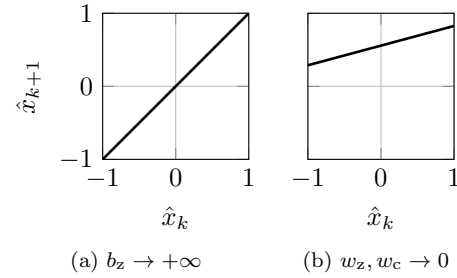


Fig. 7. Parameterizations for which GRU converges to linear activation functions.

ing different limits in the parameter space is especially important when identifying technical systems. Even for highly nonlinear processes the one step prediction surface becomes more and more linear with increasing sampling rate (Nelles, 2001). Therefore even when estimating a model for a nonlinear process, it is important that the model has the capability to represent a linear dynamical system. The fact that (5) represents a linear model for an infinite number of parameter combinations and even when approaching different limits, means, that there are huge spaces in the parameter space, which already represent a formidable solution. Representing the identity function with a one dimensional Elman Network is only possible by letting w_x and b_x approach zero, which corresponds to only a single point in the parameter space.

5. ACADEMIC EXAMPLE

Based on the idea, that representing an approximately linear system is important even when identifying a nonlinear system, the loss functions and their gradients produced by the RNN and the GRU will be examined on the following linear discrete-time autonomous system with one state:

$$x_{k+1} = -0.9 \cdot x_k \quad (13)$$

The initial state is $x_0 = 1$ and the system is simulated for $N = 100$ time steps. The RNN is equipped with one neuron in the recurrent layer without an output layer:

$$\hat{x}_{k+1} = \tanh(w_x \hat{x}_k + b_x) \quad (14)$$

Since the GRU has two gates with respectively two parameters, the parameters of the update gate f_z will be fixed to the reasonable values $w_z = 1$ and $b_z = 0$ when examining the loss function and gradient with respect to the output gates parameters. Conversely the parameters of the output gate f_c are fixed to the values $w_c = 1$ and $b_c = 0$ when examining the update gate.

$$\hat{x}_{k+1} = [\sigma(w_z \hat{x}_k + b_z) \odot \hat{x}_k + [1 - \sigma(w_z \hat{x}_k + b_z)] \cdot \tanh(w_c \hat{x}_k + b_c)] \quad (15)$$

Figures 8 and 9 show the resulting loss functions of the GRU and the RNN as well as their respective gradients.

It is apparent, that the optimization problem posed by the RNN is rather difficult to solve with gradient based optimization techniques: The local optimum lies in a narrow valley surrounded by large gradients. The loss function of the GRU on the other hand does not possess such a narrow value, compare figure 9a. The reason for this is that the update gate f_z dominates the behavior of the model for $\hat{x} > 0$ (which is where training data is

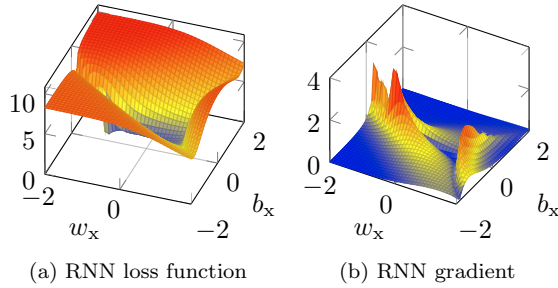


Fig. 8. RNNs loss function and magnitude of the gradient on the linear identification task

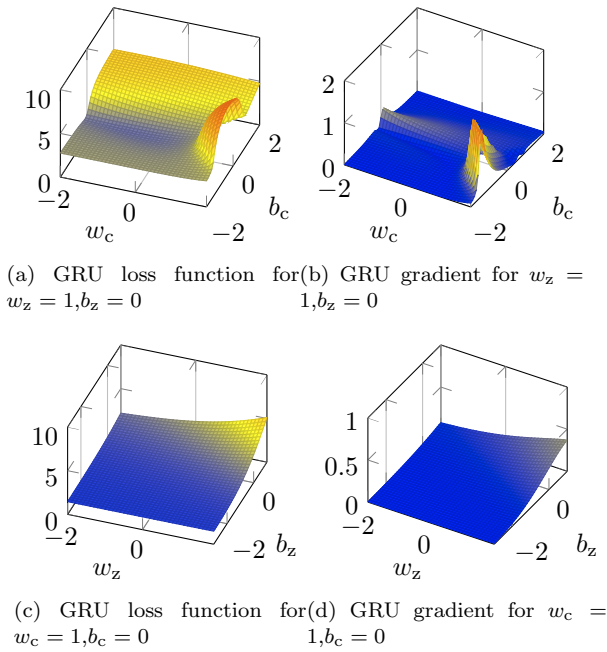


Fig. 9. GRUs loss function and magnitude of the gradient on the linear identification task

available in this case study) by multiplying the output of the output gate f_c with $1 - f_z$ and thereby diminishing its influence. The effect is a smooth loss function without large gradients.

6. CASE STUDY: ELECTRO-MECHANICAL THROTTLE

To test, whether the properties of the GRU also proof beneficial in real-life applications, it was compared to an RNN on a real nonlinear dynamical system.

6.1 The test system

The system to be identified is an industrial electro-mechanical throttle, as they are employed in combustion engines, which is operated load-free in a laboratory test stand as depicted in figure 10. The input signal u is the duty cycle for the pulse width modulator (PWM). The output signal y is the voltage of the sensor measuring the angle of the throttle plate. The throttle is actuated by a DC motor. Even though the setup is relatively simple, the modeling task is rather difficult due to the characteristics

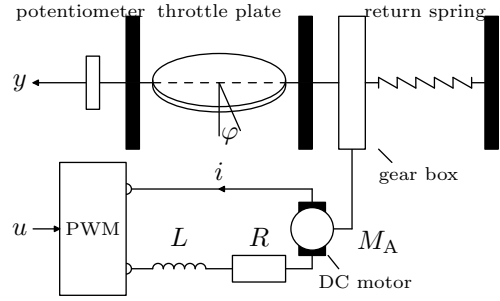


Fig. 10. Technology schematic of the electro-mechanical throttle.

of the system (Gringard and Kroll, 2016): a lower and an upper hard stop, state-dependent friction, and a nonlinear return spring.

6.2 Excitation Signals

One multisine signal and two Amplitude Modulated Pseudo-Random Binary Sequences (APRBS 1 and APRBS 2) have been used to excite the system. The multisine signal has a length of ≈ 10 s or 10^3 instances, and the APRBS signals have a length of ≈ 25 s or 2500 instances each. For the multisine signal, an upper frequency of $f_u = 7.5$ Hz has been used. For the APRBS signals, the holding time is $T_H = 0.1$ s. See (Gringard and Kroll, 2016) for more information on the test signal design.

6.3 Data preprocessing

APRBS 1 and its response signal were scaled to the interval $[-1, 1]$; all other signals were scaled accordingly. The data was then divided into training, validation and test datasets in the following way:

- Training dataset: Consists of two batches. The first batch comprises 80 % of all instances of the multisine signal and the corresponding system response. The second batch consists of 70 % of APRBS 1 and its corresponding response signal.
- Validation dataset: Consists of two batches. The first batch comprises the remaining 20 % of the multisine signal and the corresponding system response. The second batch consists of the remaining 30 % of APRBS 1 and the response signal.
- Test dataset: One batch. APRBS 2 and its corresponding system response.

This division was chosen because the multisines response signal almost exclusively covers the medium operating range while the APRBS' response signal also covers the lower and upper hard stops.

6.4 Model Architectures

The different model architectures used for identification are listed in TABLE 1. Each model architecture was initialized randomly 20 times and trained until convergence in order to examine sensitivity to initial parameters. The number of states $\dim(\mathbf{x})$ was varied from three to ten. Model architectures with an equal number of states were

designed to have the same number of hidden neurons in the output layer. Due to its many gates a GRU will have significantly more parameters than an RNN with an equal number of states. This should be considered when comparing both architectures.

Table 1: Network architectures

GRU							
1 st Layer	Gated Unit dim(\mathbf{x})	3	4	5	6	8	10
2 nd Layer	$f_g(\cdot)$ #(Neurons)	4	5	6	7	8	10
	dim(θ)	66	103	149	201	321	481
RNN							
1 st Layer	$f_h(\cdot)$ dim(\mathbf{x})	3	4	5	6	8	10
2 nd Layer	$f_g(\cdot)$ #(Neurons)	4	5	6		7	
	dim(θ)	36	55	78	105	151	205

6.5 Model Training

Each of the models in Table 1 was initialized randomly 20 times and trained for 800 epochs. Between each batch, the models initial states were set to zero. Parameters were estimated based on the training dataset using the ADAM optimizer (Kingma and Ba, 2015) with its default parameter configuration ($\alpha = 0.01, \beta_1 = 0.9, \beta_2 = 0.999$).

6.6 Results

At the end of the optimization procedure, the model with the parameter configuration, which performed best on the validation dataset, was selected and evaluated on the test dataset. The performance of the models is measured in terms of their best fit rate (BFR):

$$\text{BFR} = 100\% \cdot \max \left(1 - \frac{\|y_k - \hat{y}_k\|_2}{\|y_k - \bar{y}\|_2}, 0 \right) \quad (16)$$

Figure 11 shows the BFR of the RNN and the GRU on the test dataset. As expected, nonlinear optimization of the GRU consistently yields high performing models, while the performance of the RNNs fluctuates strongly. It should be noted, that there are cases where the RNNs performance matches or even exceeds the performance of the GRU (e.g. for $\dim(\mathbf{x}) = 10$). This proves, that the RNN is in general able to represent the test system, it seems just very unlikely to arrive at such a parameterization during the optimization process. Arguably because of the issues discussed in Section 3.

7. CONCLUSIONS & OUTLOOK

It was shown, that the gradient of the GRU's state equation w.r.t. its parameters is at most as large as but usually smaller than that of the RNN, provided the L_1 norm of all weights is smaller or equal to one. This finding contradicts the argument, that a vanishing gradient is responsible for the RNNs poor performance on various tasks. The first point made in this paper is, that the

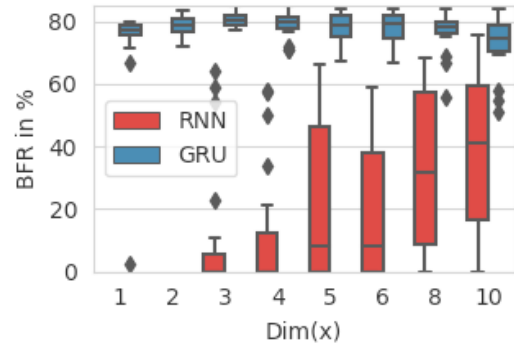


Fig. 11. Boxplot of the BFR of RNN and GRU on the test dataset. Each model was initialized 20 times and trained for 800 epochs.

smaller gradient produced by the GRU helps gradient based optimization, since small changes in the parameter space correspond to small changes in the evolution of the state, which in turn produces a smooth loss function without large gradients. The second argument made is, that the GRU's state equation converges to different functions, when certain parameters become larger. This corresponds to producing large planes in the loss function along which the solution improves steadily, rather than narrow valleys, as they are produced by the RNN. The analyses provided in this paper have yet to be generalized to state space networks with arbitrary dimensions and the whole parameter space.

REFERENCES

- Cho, K. et al. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 8, 1724–1734.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Gringard, M. and Kroll, A. (2016). On the systematic analysis of the impact of the parametrization of standard test signals. In *IEEE Symposium Series of Computational Intelligence 2016*. IEEE, Athens, Greece.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- Jordan, I.D., Sokol, P.A., and Park, I.M. (2019). Gated recurrent units viewed through the lens of continuous time dynamical systems. *arXiv preprint arXiv:1906.01005*.
- Kingma, D. and Ba, J. (2015). Adam: A method for stochastic optimization. In *3rd International Conference for Learning Representations (ICLR 2015)*.
- Nelles, O. (2001). *Nonlinear System Identification: From Classical Approaches to Neural Networks and Fuzzy Models*. Springer, Berlin Heidelberg, Germany.
- Pascanu, R., Mikolov, T., and Bengio, Y. (2012). Understanding the exploding gradient problem. *CoRR*, abs/1211.5063.
- Rehmer, A. and Kroll, A. (2019). On using gated recurrent units for nonlinear system identification. In *Preprints of the 18th European Control Conference (ECC)*, 2504–2509. IFAC, Naples, Italy.