

THESIS TITLE

MSc Thesis (*Afstudeerscriptie*)

written by

Krsto Proroković

(born March 12, 1993 in Kotor, Montenegro)

under the supervision of **Dr Germán Kruszewski** and **Dr Elia Bruni**, and submitted to the
Board of Examiners in partial fulfillment of the requirements for the degree of

MSc in Logic

at the *Universiteit van Amsterdam*.

Date of the public defense: **Members of the Thesis Committee:**

Date of the defense goes here Committee goes here



INSTITUTE FOR LOGIC, LANGUAGE AND COMPUTATION

Abstract

Abstract goes here

Acknowledgments

Contents

1	Introduction	3
2	Related Work	4
3	Background	5
3.1	Locally k-testable languages	5
3.2	Recurrent Neural Networks	5
3.2.1	Vanilla Recurrent Neural Networks	5
3.2.2	Backpropagation Through Time and Vanishing Gradient	7
3.2.3	Long-Short Term Memory Networks	7
3.2.4	Gated Recurrent Units	8
3.2.5	One-hot and Dense Encodings	8
3.2.6	Attention Mechanism	9
4	Experiments and Results	10

Chapter 1

Introduction

Chapter 2

Related Work

Learning to decide a formal language can be seen as an example of learning an algorithm from input/output examples. There are several approaches to this problem. One would be to induce a discrete program from finite set of instructions. Inductive logic programming [10] requires domain specific knowledge about the programming languages and hand-crafted heuristics to speed up the underlying combinatorial search. Levin [7] devised an asymptotically optimal method for inverting functions on given output, albeit with a large constant factor. Hutter [5] reduced the constant factor to less than 5 at the expense of introducing a large additive constant. These methods are not incremental and do not make use of machine learning. Schmidhuber [12] extended the principles of Levin’s and Hutter’s method and developed an asymptotically optimal, incremental method that learns from experience. However, the constants involved in the method are still large which limits its practical use.

The approach that we will take over here is based on training differentiable neural networks. To the best of our knowledge no one used neural networks for deciding a whole class of formal languages, rather a single language from positive and negative examples. Mozer and Das [9] used neural networks with external stack memory to parse simple context-free languages such as $a^n b^n$ and parenthesis balancing from both positive and negative examples. Wiles and Elman [14] used neural networks to learn sequences of the form $a^n b^n$ and generalise on a limited range of n . Joulin and Mikolov [6] used recurrent neural networks augmented with a trainable memory to predict the next symbol of the element of context-free language. Lastly, Zaremba and Sutskever [15] used recurrent neural networks for learning to execute simple Python programs.

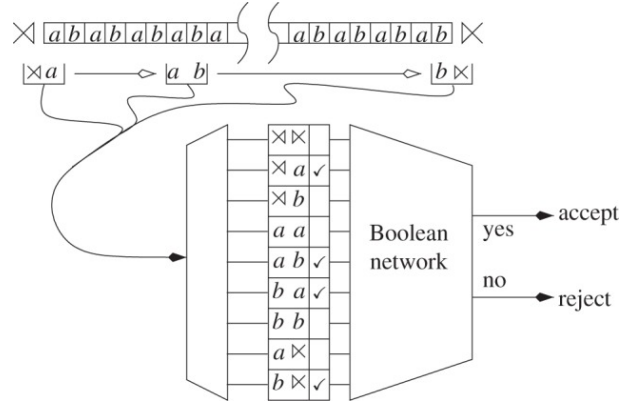
Chapter 3

Background

3.1 Locally k-testable languages

Fix an alphabet $\mathcal{A} = \{a_1, a_2, \dots, a_n\}$. An SL_k definition is just a set of blocks of k adjacent symbols (called k -factors) drawn from the alphabet. A string satisfies the description if and only if every k -factor that occurs in the string is licensed by the definition. We can see k -factors as atomic properties of strings: a string satisfies a k -factor if and only if that factor occurs somewhere in the string. Then, we can build descriptions as propositional formulas over these atoms. We will call these formulas k -expressions. A k -expression defines the set of all strings that satisfy it. A language that is defined in this way is called a locally k -testable (LT_k) language.

A scanner for an LT_k language contains a table in which it records, for every k -factor over the alphabet, whether or not that k -factor has occurred somewhere in the string. It then feeds this information into a Boolean network which implements some k -expression. When the end of the string is reached, the automaton accepts or rejects the string depending on the output of the network.



3.2 Recurrent Neural Networks

Recurrent neural networks (RNNs) [11] are parametric models of computation for processing sequences loosely inspired by biological neural networks. They found applications in handwriting recognition [reference], speech recognition [reference], machine translation [reference], etc. In this section we provide an introduction to RNNs and ... For a more detailed treatment we point the reader to [3].

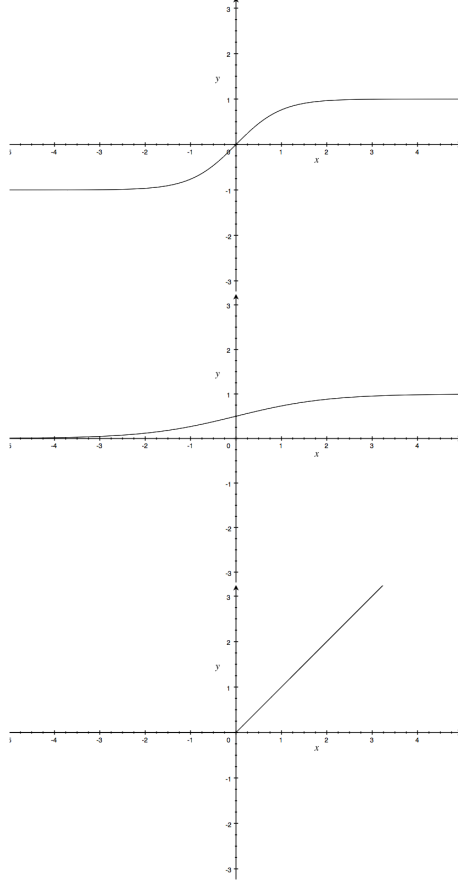
3.2.1 Vanilla Recurrent Neural Networks

We start with a simple vanilla RNN model. A vanilla RNN is given with the following parameters:

- Input to hidden connection weights $W_x \in \mathbb{R}^{d_h \times d_x}$

- Hidden to hidden connection weights $W_h \in \mathbb{R}^{d_h \times d_h}$
- Bias term $b_h \in \mathbb{R}^{d_h}$
- Activation function $\phi : \mathbb{R} \rightarrow \mathbb{R}$
- Hidden to output connection weights $W_y \in \mathbb{R}^{d_y \times d_y}$
- Bias term $b_y \in \mathbb{R}^{d_y}$
- Initial hidden state $h^{(0)} \in \mathbb{R}^{d_h}$ (usually set to zero vector)

Connection weights model the strength of synapses and activation function models neuronal firing [8]. Most commonly used activation functions are sigmoid $\sigma(x) = 1/(1 + \exp(-x))$, hyperbolic tangent $\tanh(x) = (\exp(x) - \exp(-x))/(\exp(x) + \exp(-x))$ and rectified linear unit $\text{ReLU}(x) = \max\{0, x\}$.



The computation is performed in the following way: for $t = 1, 2, \dots, T$

$$h^{(t)} = \phi \left(W_h h^{(t-1)} + W_x x^{(t)} + b_h \right)$$

$$y^{(t)} = \text{softmax} \left(W_y h^{(t)} + b_y \right)$$

where $\phi(v)_i = \phi(v_i)$ and $\text{softmax}(v)_i = \exp(v_i) / \sum_j \exp(v_j)$.

RNN figure goes here.

Note that components of $y^{(t)}$ sum up to one, so we can interpret it as a probability distribution. Here we will be mostly interested in sequence classification, i.e. we will only care about $y^{(T)}$.

Example Consider a vanilla RNN given with the following parameters

$$W_x = \begin{bmatrix} 3 & -1 & 0 \\ 0 & 0 & 2 \end{bmatrix} \quad W_h = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad b_h = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad W_y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad b_y = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

and $\phi(z) = z$. It is easy to see that $y_1^{(T)} < y_2^{(T)}$ if and only if $\sum_t (3x_1^{(t)} - x_2^{(t)}) < \sum_t 2x_3^{(t)}$. Note that in this case $h^{(t)}$ accumulates the results of the computation on the input encountered so far.

While the previous example displays a rather simple function, vanilla RNNs are universal in the sense that any function computable by a Turing machine can be computed by such a network [reference Siegelmann and Sontag].

We can also stack several RNN cells in order to exploit compositionality and obtain more powerful RNNs. Such RNNs are called multilayer or deep RNNs. The computation performed by an L layer vanilla RNN is given with

$$\begin{aligned} \mathbf{h}_1^{(t)} &= \phi \left(\mathbf{W}^{\mathbf{h}_1 \mathbf{h}_1} \mathbf{h}_1^{(t-1)} + \mathbf{W}^{\mathbf{x} \mathbf{h}_1} \mathbf{x}^{(t)} + \mathbf{b}^{\mathbf{h}_1} \right) \\ \mathbf{h}_2^{(t)} &= \phi \left(\mathbf{W}^{\mathbf{h}_2 \mathbf{h}_2} \mathbf{h}_2^{(t-1)} + \mathbf{W}^{\mathbf{x} \mathbf{h}_2} \mathbf{x}^{(t)} + \mathbf{b}^{\mathbf{h}_2} \right) \\ &\vdots \\ \mathbf{h}_L^{(t)} &= \phi \left(\mathbf{W}^{\mathbf{h}_L \mathbf{h}_L} \mathbf{h}_L^{(t-1)} + \mathbf{W}^{\mathbf{x} \mathbf{h}_L} \mathbf{x}^{(t)} + \mathbf{b}^{\mathbf{h}_L} \right) \\ \mathbf{y}^{(t)} &= \text{softmax} \left(\mathbf{W}^{\mathbf{h}_L \mathbf{y}} \mathbf{h}_L^{(t-1)} + \mathbf{b}^{\mathbf{y}} \right) \end{aligned}$$

Observe that we will need L initial states $\mathbf{h}_1^{(0)}, \mathbf{h}_2^{(0)}, \dots, \mathbf{h}_L^{(0)}$.

3.2.2 Backpropagation Through Time and Vanishing Gradient

In practice we need to estimate the parameters of RNN from the data. The way this is done is by minimising the average of a certain loss function which measures how far are the outputs of the network from actual data. Suppose that we are given a dataset $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$. The quantity that we are supposed to minimise is

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i)$$

where θ are parameters of the network and $\hat{\mathbf{y}}_i$ is the output of the network on input \mathbf{x}_i . The loss function that we will use here is the cross entropy loss given with

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i.$$

In order to minimise the loss function we need to calculate its gradient with respect to the parameters of the network. This is achieved by Backpropagation Through Time algorithm [13].

Unfortunately, for vanilla RNNs this does not solve the optimisation problem. The reason is that is that the gradient of the loss function may vanish or explode over many time steps. This problem has been studied independently by separate researchers [Hochreiter][Bengio]. For survey see [paper by Pascanu and Bengio].

The problem of exploding gradient can be solved by gradient clipping: if the norm of the gradient exceeds a certain threshold, then normalise it and multiply it by that threshold. However, the vanishing gradient remains a problem for vanilla RNNs which makes their optimisation difficult in practice.

3.2.3 Long-Short Term Memory Networks

One way to deal with vanishing gradient is to use gated RNNs whose connection weights may change at each time step. Multiplicative gates allow RNNs to control the way information is stored and retrieved. Thus, they are able to create paths through time that have derivatives that neither vanish nor explode. Gates in RNNs can be seen as differentiable version of logic gates in digital circuits.

Long short term memory (LSTM) networks [4] use additional memory cell $c^{(t)}$ apart from $h^{(t)}$. Here we will be using standard LSTM architecture introduced in [2] which uses four gates: forget

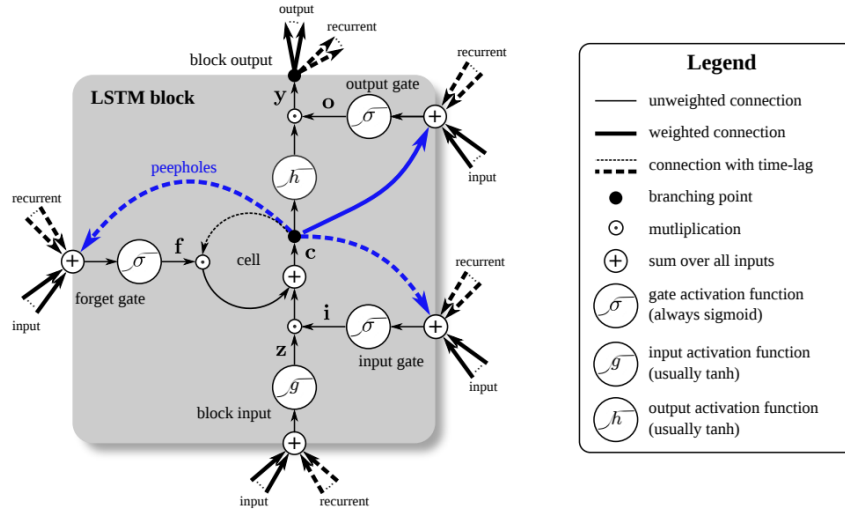
gate $f^{(t)}$, input gate $i^{(t)}$, output gate $o^{(t)}$ and input modulation gate $g^{(t)}$. Its cell is described by the following set of equations

$$\begin{bmatrix} f^{(t)} \\ i^{(t)} \\ o^{(t)} \\ g^{(t)} \end{bmatrix} = W_x x^{(t)} + W_h h^{(t)} + b$$

$$c^{(t)} = \sigma(f^{(t)}) \odot c^{(t-1)} + \sigma(i^{(t)}) \odot \tanh(g^{(t)})$$

$$h^{(t)} = \sigma(o^{(t)}) \odot \tanh(c^{(t)})$$

where \odot is component-wise or Hadamard product.



Just like vanilla RNN cells, LSTM cells can be stacked in order to obtain more powerful multilayer LSTM networks.

LSTM networks have been successfully applied in many contexts such as handwriting recognition [reference] and generation, language modeling [reference], machine translation [reference], image captioning [reference], parsing [reference], etc.

3.2.4 Gated Recurrent Units

Another gated recurrent architecture frequently used in practice are Gated Recurrent Units (GRU) [1]. Unlike LSTM, GRU gets rid of memory cell $c^{(t)}$ by combining input and forget gates into a single one. The GRU cell is described by the following set of equations

$$r^{(t)} = \sigma(W_{xr}x_t + W_{hr}h^{(t-1)} + b_r)$$

$$z^{(t)} = \sigma(W_{xz}x_t + W_{hz}h^{(t-1)} + b_z)$$

$$n^{(t)} = \tanh(W_{xn}x_t + r_t(W_{hn}h^{(t-1)} + b_n))$$

$$h^{(t)} = (1 - z_t)n_t + z_th_{t-1}$$

GRU figure goes here.

Again, GRU cells can be stacked, yielding multilayer GRU networks.

3.2.5 One-hot and Dense Encodings

Note that RNNs work exclusively with numerical data, i.e. real numbers. Suppose that we are given an alphabet $\mathcal{A} = \{a_1, a_2, \dots, a_k\}$. One way to encode it would be to map each a_i to i . However, this is not a good encoding as the network would need to learn that there is no ordering between alphabet

elements. A better way to encode the alphabet that avoids this problem is to assign a one-hot vector to every element of the alphabet, i.e.

$$a_1 \mapsto \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad a_2 \mapsto \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \dots \quad a_k \mapsto \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

One problem with this encoding is that k may be huge while one-hot vectors remain sparse. For example, in natural language processing k would be the number of words which is usually several tens of thousands. This problem is solved by using embedding layer which learns to map alphabet elements to dense, low-dimensional encodings.

3.2.6 Attention Mechanism

When processing the input, RNNs store all the memory in the hidden state. However, it may be hard to compress potentially long input in a single context vector. This could definitely be an issue if the hidden state is passed as input to another network which is common in sequence to sequence learning. For example, the family of models used for machine translation has one RNN called encoder to process the input in one language and pass it to the other RNN called decoder which is supposed to output a sequence in a different language.

Attention mechanism allows us to obtain a distribution on part of the input provided so far and focus on certain parts of it in a way we humans do. The attentional hidden state is defined as

$$\tilde{h}^{(t)} = \tanh\left(W_c[c^{(t)}; h^{(t)}]\right)$$

where $c^{(t)}$ is the context vector and $[:]$ denotes concatenation. We will be using an attention model in which the alignment vector $a^{(t)}$ is obtained by comparing the current hidden state $h^{(t)}$ with each source hidden state $\bar{h}^{(s)}$

$$a^{(t)} = \text{align}(h^{(t)}, \bar{h}^{(s)}) = \text{softmax}(\text{score}(h^{(t)}, \bar{h}^{(s)})) .$$

Attention figure goes here.

Chapter 4

Experiments and Results

Bibliography

- [1] Kyunghyun Cho et al. “Learning phrase representations using RNN encoder-decoder for statistical machine translation”. In: *arXiv preprint arXiv:1406.1078* (2014).
- [2] Felix A. Gers, Jürgen Schmidhuber, and Fred Cummins. “Learning to forget: Continual prediction with LSTM”. In: (1999).
- [3] Alex Graves. “Supervised sequence labelling”. In: *Supervised sequence labelling with recurrent neural networks*. Springer, 2012, pp. 5–13.
- [4] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Computation* 9.8 (1997), pp. 1735–1780.
- [5] Marcus Hutter. “The fastest and shortest algorithm for all well-defined problems”. In: *International Journal of Foundations of Computer Science* 13.03 (2002), pp. 431–443.
- [6] Armand Joulin and Tomas Mikolov. “Inferring algorithmic patterns with stack-augmented recurrent nets”. In: *Advances in neural information processing systems*. 2015, pp. 190–198.
- [7] Leonid Anatolevich Levin. “Universal sequential search problems”. In: *Problemy Peredachi Informatsii* 9.3 (1973), pp. 115–116.
- [8] Rodolfo Llinas. “Neuron”. In: *Scholarpedia* 3.8 (2008). revision #91570, p. 1490. DOI: [10.4249/scholarpedia.1490](https://doi.org/10.4249/scholarpedia.1490).
- [9] Michael C Mozer and Sreerupa Das. “A connectionist symbol manipulator that discovers the structure of context-free languages”. In: *Advances in neural information processing systems*. 1993, pp. 863–870.
- [10] Stephen Muggleton. “Inductive logic programming”. In: *New generation computing* 8.4 (1991), pp. 295–318.
- [11] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. *Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [12] Jürgen Schmidhuber. “Optimal ordered problem solver”. In: *Machine Learning* 54.3 (2004), pp. 211–254.
- [13] Paul J. Werbos. “Backpropagation through time: what it does and how to do it”. In: *Proceedings of the IEEE* 78.10 (1990), pp. 1550–1560.
- [14] Janet Wiles and Jeff Elman. “Learning to count without a counter: A case study of dynamics and activation landscapes in recurrent networks”. In:
- [15] Wojciech Zaremba and Ilya Sutskever. “Learning to execute”. In: *arXiv preprint arXiv:1410.4615* (2014).