

VICTORIA UNIVERSITY OF WELLINGTON
Te Whare Wānanga o te Ūpoko o te Ika a Māui



School of Engineering and Computer Science
Te Kura Mātai Pūkaha, Pūrorohiko

PO Box 600
Wellington
New Zealand

Tel: +64 4 463 5341
Fax: +64 4 463 5045
Internet: office@ecs.vuw.ac.nz

some kind of rnn/tensor mess

Paul Francis Cunninghame Mathews

Supervisors: Marcus Frean and David Balduzzi

Submitted in partial fulfilment of the requirements for
Bachelor of Science with Honours in Computer Science.

Abstract

A short description of the project goes here.

Acknowledgments

Any acknowledgments should go in here, between the title page and the table of contents. The acknowledgments do not form a proper chapter, and so don't get a number or appear in the table of contents.

WHO??

Andrej Karpathy, for generously “lending” a figure.

Contents

1	Introduction	1
2	Background and Related Work	3
2.1	Feed-Forward Neural Networks	3
2.1.1	Training	4
2.2	Recurrent Neural Networks	5
2.2.1	Classical Formulation	5
2.2.2	Training	5
2.2.3	Issues with Vanilla RNNs	6
2.3	Popular Alternatives: LSTM and GRU	7
2.4	Other Ideas	8
2.4.1	Enhancing Training	8
2.4.2	Memory	10
2.5	Tensors in Neural Networks	11
3	Three-way Tensors and Bilinear Products	13
3.1	Definitions	13
3.2	Interpreting Bilinear Products	15
3.2.1	Stacked bilinear forms	15
3.2.2	Choosing a matrix	16
3.2.3	Tensor as an independent basis	17
3.2.4	Operation on the outer product	17
3.3	Tensor Decompositions	21
3.3.1	CANDECOMP/PARAFAC	22
3.3.2	Tensor Train	24
3.3.3	Comparison	25
3.4	Learning decompositions by gradient descent	27
3.4.1	Gradients	27
3.4.2	Random Bilinear Products	28
3.4.3	Learning to multiply	29
3.4.4	XOR	31
3.4.5	Separating Style and Content	33
4	Proposed Architectures	37
4.1	Incorporating tensors for expressivity	37
4.1.1	Biases	37
4.1.2	Generalised Multiplicative RNN	38
4.2	Gates and Long Time Dependencies	38
4.2.1	Naive Addition	38
4.2.2	Gates	39

4.3	Gated Tensor RNNs	45
4.3.1	Choices	45
4.3.2	Tensor Gate Unit	47
5	Evaluation of architectures	49
5.1	Synthetic Tasks	49
5.1.1	Addition	49
5.1.2	Variable Binding	52
5.1.3	Sequential MNIST	55
5.2	Real-world Data	60
5.2.1	Polyphonic Music	60
5.2.2	PTB	63
5.2.3	War and Peace	65
6	Conclusions	67
A	Additional Proofs	69
B	Algorithms for Synthetic Tasks	71
B.1	Addition	71
B.2	Variable binding	71

Figures

3.1	Example tensor network diagrams.	14
3.2	Various products expressed as Tensor Network Diagrams.	15
3.3	One way of representing a bilinear product in the CP-decomposition.	23
3.4	Diagrams of the Tensor Train decomposition.	25
3.5	Results for learning direct approximations of a CP-rank 100 tensor.	30
3.6	Results for learning direct approximations of a tt-rank 16 tensor.	30
3.7	Training error for various rank decompositions on the element-wise multiplication task.	32
3.8	Training error for various rank decompositions on the permuted element-wise multiplication task.	32
3.9	Training error for two decompositions learning element-wise multiplication with varying amounts of structure.	33
3.10	XOR results	34
3.11	Learned style and content representations.	36
3.12	Example of images generated from unseen style and content pairs, three letters from three different fonts.	36
3.13	A difficult example.	36
4.1	Forget gate window shapes	41
4.2	Saturations of 3-layer LSTM gates	41
4.3	Convex gate shapes	44
4.4	Convex gate choosing several past candidates	44
5.1	Results for addition task, standard sequence lengths.	51
5.2	Results for addition task, very long sequences.	52
5.3	Example instances of variable binding	53
5.4	Example successful outputs for variable binding	53
5.5	Example of baseline for variable binding	53
5.6	Variable binding results, one pattern	56
5.7	Variable binding results, two patterns	56
5.8	Variable binding results, three patterns	56
5.9	Permuted MNIST results	59
5.10	TGU PTB results by rank	64
5.11	Training and validation curves on War and Peace. Dashed lines indicate training loss, solid lines are the loss on the validation set. 'Full' refers to rank equal to the number of hidden units, 'one' simply means rank 1.	66

Chapter 1

Introduction

Recurrent Neural Networks (RNNs) are powerful machine learning models for learning mappings between pairs of sequences. Such tasks are typically highly challenging and include language modelling, machine translation, speech recognition and audio classification. They have achieved remarkable success even at immense scale [103], but the most successful architectures are opaque, poorly understood and seemingly redundant while simple clear models tend to have serious issues learning complex sequence tasks [6]. Understanding how RNNs are able to solve complex tasks and designing new, simpler architectures that are well-understood from the outset is therefore of high importance.

RNNs make computations at each time step of the input sequence based on the current input value and a fixed size vector of hidden states. At each step a new state is produced and passed forward. This allows the network to model complicated temporal dynamics. The most successful architectures have two key components to this computation: additive state updates and multiplicative gating. The former allows easy integration of new information at each time step while the latter provides the ability to selectively ignore elements of the state or input.

Multiplicative gating involves computing element-wise products of vectors. We can generalise this as computing a product involving a three-way tensor of weights, termed a bilinear product. Looking to understand and generalise the structures prevalent in existing RNNs we investigate in detail the properties of the bilinear product. This leads to the remarkable result, expressed in theorem 1, that these bilinear products operate on the exclusive-or of features. The implications of this are clear and strong – bilinear products operate implicitly at a higher level than the standard linear operations composed to form neural networks.

The downside of these tensor products is that using them explicitly requires storing the tensor which can be very large. We therefore look to the multi-linear algebra literature to investigate methods of tensor decomposition. Finding a useful decomposition would allow us to incorporate this highly expressive tensor product in a neural network with a feasible number of parameters. We find that the simplest decomposition, CANDECOMP/PARAFAC [9, 35], fulfils our criteria both theoretically and experimentally.

Equipped with this tensor decomposition we now propose two new classes of RNNs: the Generalised Multiplicative RNN (GMR) and the Tensor Gate Unit (TGU). The GMR is a very simple way to incorporate a tensor into a neural network and generalises a number of recently proposed architectures [91, 104]. The TGU includes an additive state update and standard multiplicative gates in addition to the bilinear product. The architecture is based on the intuition. Facilitating this suggests a distinct architecture with a unique and clear motivation.

These architectures are then evaluated empirically. We consider standard synthetic tasks and achieve excellent results including successfully training a TGU to capture time

dependencies up to 10,000 steps, more than 10 times the best previously reported. None of the synthetic tasks common in the literature adequately exercise the ability of the model to store and retrieve arbitrary patterns. As this is an essential sub-task for solving many complex problems we propose a novel synthetic form of *variable binding* which tests the ability of the models to use their hidden units in this fashion.

Finally, we demonstrate that our proposed RNNs are highly competitive with the best existing techniques on real world data. In doing so, we investigate the hypothesis that controlling the rank of the tensor decomposition can provide an important regularisation effect, helping greatly to reduce overfitting. These results justify the intuitions and theoretical analysis behind the proposed architectures as well as indicating they provide a useful class of novel RNNs with clear performance benefits.

Chapter 2

Background and Related Work

2.1 Feed-Forward Neural Networks

Feed-forward neural networks are a class of parameterised function approximators which consist of artificial neurons arranged in layers. They arise from generalisations of the perceptron [81], indeed simple feed-forward networks are often termed ‘multi-layer perceptrons’ (MLPs). MLPs solve the key limitations of the early perceptron [68], but it was not until 1986 when a reliable and general purpose procedure (*back-propagation*) was proposed [82] to train them.

Conceptually it is possible to consider MLPs as a directed acyclic graph, where each node consists of a ‘neuron’ or ‘unit’ with the direction of the edges indicating the way information flows from the input of the function being modelled to the output. For an MLP this graph is divided into layers. The first layer represents the input (with one unit per coordinate) while subsequent layers represent individual neurons. In a standard fully-connected network, every node is connected to all nodes in the subsequent layer. When an input is presented to the network, each neuron computes an activation by a weighted sum over its inputs followed by the application of some non-linear function. Specifically a unit that receives inputs (x_1, x_2, \dots, x_n) will compute its activation h as:

$$h = f \left(\sum_{i=1}^n x_i w_i \right) \quad (2.1)$$

where $(w_1, w_2, \dots, w_n) + b$ are the weights the neuron gives its inputs and b is a bias, commonly added to enhance the expressive power. $f : \mathbb{R} \rightarrow \mathbb{R}$ is some nonlinear function, common examples would be a threshold, the logistic sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ or a linear rectifier $\rho(x) = \max(0, x)$ [69].

While it is convenient to consider the network in this manner, it is computationally and notationally more efficient to consider to express many neurons at once, using vectors and matrices. If we consider the input to a network as a vector of features \mathbf{x} , we can express eq. (2.1) with a matrix multiplication as

$$h = f(\mathbf{w}^T \mathbf{x} + b)$$

where \mathbf{w} is the weights from eq. (2.1) as a vector. We can generalise this to compute an entire layer of activations at once. Let \mathbf{W} be a matrix whose rows contain the weight vectors for each neuron in a given layer and \mathbf{b} a vector containing all of their biases. We can then compute their activations at once with

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$

provided we ensure f is applied elementwise to the result.

This gives us a single layer, a classic MLP will perform this process twice, although the last layer is often linear:

$$\hat{\mathbf{y}} = \mathbf{W}^{(2)} f(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}.$$

This can be thought of as proceeding in stages – first compute \mathbf{h} as before and then using another set of weights compute a new transformation of it to get the final output.

Although it can not increase the class of functions a neural network can approximate, it is often beneficial in practice to increase the depth as this can allow to express more complex relationships with only a moderate increase in parameters to be learned [96]. This is done by recursively applying the above transform with a fresh set of weights and biases. Omitting the biases for clarity this gives us the following for a network with $l - 1$ hidden layers:

$$\hat{\mathbf{y}} = \mathbf{W}^{(l)} f(\mathbf{W}^{(l-1)} f(\dots f(\mathbf{W}^{(1)} \mathbf{x}))).$$

Naively adding depth in this fashion can cause the network to become very challenging to train [27, 84]. Several methods have been proposed which successfully alleviate this, such as layer-wise pre-training [41, 99] and more recently residual connections which change the structure of the graph to allow better flow of information [36].

2.1.1 Training

These networks are typically trained in a supervised fashion. Denote the parameters of the network $\theta = \{\mathbf{W}^{(l)}, \dots, \mathbf{W}^{(1)}, \mathbf{b}^{(l)}, \dots, \mathbf{b}^{(1)}\}$ and the output of the network for a given input \mathbf{x} and set of parameters as $\hat{\mathbf{y}} = \mathcal{F}(\mathbf{x}; \theta)$. Suppose we have a set of m data items $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$ where each item is an input-output pair. Let \mathcal{X} be the set of possible \mathbf{x} values and \mathcal{Y} the set of possible \mathbf{y} values. We wish to learn a mapping which will produce each \mathbf{y}_i from the given \mathbf{x}_i . To do this we define a loss function $\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ which maps pairs of valid outputs to a real number, such that $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = 0$ implies $\mathbf{y} = \hat{\mathbf{y}}$. Typical loss functions are the squared error

$$\mathcal{L}_{SE}(\mathbf{y}, \hat{\mathbf{y}}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \sum_d (y_d - \hat{y}_d)^2$$

or the cross entropy

$$\mathcal{L}_{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_d y_d \log(\hat{y}_d),$$

ensuring some safeguards are put in place to avoid attempting to compute $\log 0$. The cross entropy is typically preferred when the targets can be interpreted as probabilities.¹ This requires converting the output of the network to a probability, typically done with either the sigmoid function mentioned above or the softmax $s_i(\mathbf{x}) = e^{x_i} / \sum_j e^{x_j}$.

The best set of parameters for the network, θ^* , is then the set which minimises the loss across the data items. Training the network amounts to solving the following optimisation problem:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^m \mathcal{L}(\mathbf{y}_i, \mathcal{F}(\mathbf{x}_i; \theta)).$$

¹For example, discrete class labels – we can convert them to a one-of- n (often termed “one-hot”) taking discrete labels between 1 and n to n dimensional vectors with 0 in all positions bar one, which has its element set to 1. We can then view these vectors as a (very certain) probability distribution over classes. In this case the cross entropy corresponds exactly to the negative log-likelihood of the correct class.

For neural networks this is typically done using variations on gradient descent. In its simplest form, each weight is updated at each iteration as follows:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \cdot \nabla_{\mathbf{W}} \mathcal{L}$$

where η is a hyper parameter adjusting the step size (called the *learning rate*) and $\nabla_{\mathbf{W}} \mathcal{L}$ is the gradient of the loss with respect to the \mathbf{W} . Typically the gradient is evaluated over a small subset or “mini-batch” of the data per iteration which is akin to using a noisy estimate of the true gradient. The gradients of deep networks can be calculated using *back-propagation*, which amounts to the chain rule of calculus and the observation that neural networks consist of many composed functions [82]. In practice this is assisted greatly by software packages which can determine the derivatives of a given function using automatic or symbolic differentiation such as Tensorflow [1] or Theano [98].

2.2 Recurrent Neural Networks

Recurrent Neural Networks (RNNs) of the form considered here generalise feed-forward networks to address problems in which we wish to map a *sequence* of inputs $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ to a sequence of outputs $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T)$. They have been applied successfully to a wide range of tasks which can be framed in this way including statistical language modelling [66] (including machine translation [11, 103]), speech recognition [31], polyphonic music modelling [8], music classification [13], image generation [34], automatic captioning [101, 105] and more.

2.2.1 Classical Formulation

An RNN is able to maintain context over a sequence by transferring its hidden state from one time step to the next. We refer to the vector of states at time t as \mathbf{h}_t . The classic RNN (often termed “vanilla”) originally proposed in [24] computes its hidden states with the following recurrence:

$$\mathbf{h}_t = f(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b}) \quad (2.2)$$

where $f(\cdot)$ is some elementwise non-linearity, often the hyperbolic tangent: $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$. Final outputs are then computed from these states in the same fashion as the feed-forward networks above. This remains the same for the extended architectures considered below, the key difference is the manner in which the states are produced.

Equation (2.2) bears a striking resemblance to the building block of a feed-forward network. The key difference is the (square) matrix \mathbf{W} which contains weights controlling how the previous state affects the computation of the new activations.

2.2.2 Training

We can train this (or any of the variants we will see subsequently) using back-propagation. Often termed “Back Propagation Through Time” [102] which requires using the chain rule to determine the gradients of the loss with respect to the network parameters in the same manner as for feedforward networks.

To understand what is required to perform this, consider a loss function for the whole sequence of the form

$$\mathcal{L}(\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_T, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T) = \sum_{i=1}^T \mathcal{L}_i(\hat{\mathbf{y}}_i, \mathbf{y}_i),$$

which is a sum of the loss accrued at each time step. This captures all common cases including sequence classification or regression, as the \mathcal{L}_i may simply return 0 for all but the last time step. To find gradients of the loss with respect to the parameters which generate the hidden states, we must first find the gradient of the loss with respect to the hidden states themselves. Choosing a hidden state i somewhere in the sequence we have:

$$\nabla_{\mathbf{h}_i} \mathcal{L} = \sum_{j=i}^t \nabla_{\mathbf{h}_i} \mathcal{L}_j$$

from the definition of the loss and the fact that a hidden state may affect all future losses. To determine each $\nabla_{\mathbf{h}_i} \mathcal{L}_j$ (noting $j \geq i$), we apply the chain rule, to back-propagate the error from time j to time i . This is the step from which the algorithm derives its name, and simply requires multiplying through adjacent time steps. Let $\mathbf{z}_k = \mathbf{W}\mathbf{h}_{k-1} + \mathbf{U}\mathbf{x}_k + \mathbf{b}$ be the pre-activation of the hidden states. Then

$$\begin{aligned} \nabla_{\mathbf{h}_i} \mathcal{L}_j &= \left(\prod_{k=i+1}^j \nabla_{\mathbf{h}_{k-1}} \mathbf{h}_k \right) (\nabla_{\mathbf{h}_j} \mathcal{L}_j) \\ &= \left(\prod_{k=i+1}^j \nabla_{\mathbf{z}_k} f \cdot \mathbf{W} \right) (\nabla_{\mathbf{h}_j} \mathcal{L}_j). \end{aligned} \quad (2.3)$$

This has two key components: $\nabla_{\mathbf{h}_j} \mathcal{L}_j$ quantifies the degree to which the hidden states at time j affect the loss (computing this will most likely require further back-propagation through one or more output layers) while the second term in equation (2.3) measures how much the hidden state at time i affects the hidden state at time j .

We can now derive an update rule for the parameters by observing²

$$\begin{aligned} \nabla_{\mathbf{W}} \mathcal{L} &= \sum_{i=1}^T \nabla_{\mathbf{W}} \mathcal{L}_i \\ &= \sum_{i=1}^T \sum_{j=1}^i (\nabla_{\mathbf{W}} \mathbf{h}_j) (\nabla_{\mathbf{h}_j} \mathcal{L}_i) \end{aligned} \quad (2.4)$$

and applying the above. For the input matrix \mathbf{U} and the bias the process is the same.

2.2.3 Issues with Vanilla RNNs

Vanishing and Exploding Gradients

Equation (2.3) reveals a key pathology of the vanilla RNN – vanishing gradients. This occurs when the gradient of the loss vanishes to a negligibly small value as we propagate it backward in time, leading to a negligible update to the weights. $\nabla_{\mathbf{h}_j} \mathcal{L}_j$ (a column vector) is multiplied by a long product of matrices, alternating between $\nabla_{\mathbf{z}_k} f$ and \mathbf{W} . If we assume for illustrative purposes that $f(\cdot)$ is the identity function (so we have a linear network), then the loss vector is multiplied by \mathbf{W} taken to the $(j - i)$ -th power. If the largest eigenvalue of \mathbf{W} is large, then this will cause the gradient to eventually explode. If the largest eigenvalue is small, then the gradient will vanish. This issue was first presented in 1994 [6] where it was shown to be

²For clarity, when we consider the gradient of anything with respect to a matrix, such as $\nabla_{\mathbf{W}} \mathcal{L}$ we implicitly vectorise the matrix. More precise notation would be $\nabla_{\text{vec}(\mathbf{W})} \mathcal{L}$, but this adds extra distraction into already unfortunately cluttered equations. For a definition of the vectorisation operation see section 3.2. This means $\nabla_{\mathbf{W}} \mathcal{L}$ has shape $m \times 1$, where m is the number of hidden states. Correspondingly, $\nabla_{\mathbf{W}} \mathbf{h}_j$ has shape $m \times n$, so the shapes of the matrix multiplication in eq. (2.4) are appropriately defined.

unavoidable if the network needs to both store information and be robust to noise. However, this has not precluded further research into alleviating the main symptoms, for a thorough treatment including necessary conditions for vanishing and the complementary exploding problem, see [77].

In the non-linear case, this remains a serious issue. While exploding gradients are often mitigated by using a *saturating* non-linearity so that the gradient tends to zero as the hidden states grow, this only exacerbates the vanishing problem.

The key symptom of vanishing gradients in RNNs is that it makes it much harder to learn to store information for long time periods [6]. This makes RNNs often struggle to solve simple-seeming tasks in which the solution requires remembering an input for many time steps.

Butterfly Effect

A second issue when training RNNs can be illustrated by viewing them as iterated non-linear dynamical systems and thus susceptible to the “butterfly effect”: seemingly negligible changes in initial conditions can lead to catastrophic changes after a number of iterations [58]. In RNNs this manifests as near-discontinuity of the loss surface [77] as a change (for example, to a weight during back-propagation) which may even reduce the loss for a short period can cause instabilities further on which lead to steep increases in loss. This problem is not as well studied as vanishing gradients although some partial solutions exist such as clipping the norm of gradients [77] or using a regulariser to encourage gradual changes in hidden state [54].

2.3 Popular Alternatives: LSTM and GRU

To address these fundamental problems a number of alternate architectures have been proposed. Here we will outline two popular variants: the Long Short Term Memory (LSTM) and the Gated Recurrent Unit (GRU). Both of these belong to a class of *gated* RNNs, which have a markedly different method of computing a new state.

The LSTM was proposed to alleviate the vanishing gradient problem. It uses a number of gates to control the flow of information during the computation of new states. Although a large number of variants exist, the standard LSTM we will consider here has the form: [44, 29]

$$\begin{aligned} \mathbf{h}_t &= \mathbf{o}_t \odot \tau(\mathbf{c}_t) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \\ \mathbf{g}_t &= \tau(\mathbf{W}_g \mathbf{c}_{t-1} + \mathbf{U}_g \mathbf{x}_t + \mathbf{b}_g) \\ \mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{c}_{t-1} + \mathbf{U}_o \mathbf{x}_t + \mathbf{b}_o) \\ \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{c}_{t-1} + \mathbf{U}_f \mathbf{x}_t + \mathbf{b}_f) \\ \mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{c}_{t-1} + \mathbf{U}_i \mathbf{x}_t + \mathbf{b}_i) \end{aligned}$$

where $\tau(\cdot)$ refers to the elementwise tanh, $\sigma(\cdot)$ the elementwise logistic sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ and $\cdot \odot \cdot$ is used to denote elementwise multiplication between vectors. These equations can be hard to take at face value – the key elements are $\mathbf{i}_t, \mathbf{o}_t, \mathbf{f}_t$, termed the *input*, *output* and *forget* gates respectively, are computed in the same fashion as the activations of a vanilla neural network but use sigmoid activation function which varies smoothly between zero and one. Combining this with elementwise multiplication has the eponymous gating effect, attenuating the contributions of other components.

The output gate is fairly straightforward, it simply allows the network to prevent its hidden state from being exposed. The forget and input gates have a more difficult role to characterise. Taken together, these control the acceptance or rejection of new information by modulating the amount by which the new candidate state \mathbf{g}_t is accepted into the hidden state \mathbf{c}_t .

A closely related architecture proposed much more recently by Cho et al. [11] is the GRU, which computes its state as follows:

$$\begin{aligned}\mathbf{h}_t &= \mathbf{f}_t \odot \mathbf{h}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{z}_t \\ \mathbf{z}_t &= \tau(\mathbf{W}_z(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{U}_z \mathbf{x}_t + \mathbf{b}_z) \\ \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{U}_f \mathbf{x}_t + \mathbf{b}_f) \\ \mathbf{r}_t &= \sigma(\mathbf{W}_r \mathbf{h}_{t-1} + \mathbf{U}_r \mathbf{x}_t + \mathbf{b}_r).\end{aligned}$$

This is a slightly simpler form than the LSTM although the alterations go beyond simply removing the output gate. Notably, the forget gate now controls both parts of the state update. Further, in the computation of \mathbf{z}_t there is a departure from the vanilla RNN-style building block that makes up all of the LSTM's operations. This is interesting as a half of the model's parameters, and therefore a large part of its computational power, is dedicated towards computing state updates. However, the mechanism of their computation places significant emphasis on using temporally local information to do so – the *reset* gate \mathbf{r}_t provides the model with the ability to ignore parts of its state.

The key shared component of these architectures is an *additive* state update. Another way of phrasing this is that while the vanilla RNN attempts to learn an opaque function $\mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1})$, these gated architectures instead learn a *residual* mapping $\mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}) + \mathbf{h}_{t-1}$. By itself, this would alleviate the main cause of vanishing gradients [48, 44]. This is equivalent to adding a skip connection, allowing the state to skip a time step and is directly analogous to the residual connections now commonly used to address the vanishing gradient problem in very deep feed-forward networks [36, 21, 93]. Unfortunately the presence of the gate complicates this perspective – this will be analysed in detail in chapter 4.

2.4 Other Ideas

While the LSTM and the GRU are by far the most popular RNN architectures, many other solutions have been proposed to the problems that plague the Vanilla RNN. We divide these into two categories: training and memory augmentation. Training solutions focus on adjusting the structure of the RNN or the training algorithm to avoid vanishing or exploding gradients and hence allow it to learn to solve difficult tasks. Memory augmentation approaches view an RNN as a device which interacts with a memory (in the form of its hidden states) and either augment the memory with additional structures or impose specific access patterns to encourage more appropriate use of this memory.

2.4.1 Enhancing Training

Methods of addressing issues when training RNNs can be further divided into two categories: architectural and algorithmic. Architectural solutions, such as the LSTM and GRU, seek to alter the way in which hidden states are calculated. Algorithmic solutions look to change the training process itself to better avoid or handle the various pathologies. These are complementary lines of enquiry – advanced architectures can still benefit from more sophisticated training techniques.

Architectural Solutions

An early attempt to solve the vanishing gradient problem involves feeding inputs to different units at different rates. The primary motivation for this seems to be to force the network to attend to longer time-scales by passing some parts of its information at a slower rate, although the authors note an intuition that more abstract representations of the sequence will change slower and hence not require re-calculation at every step [39]. This idea has been revisited recently in the form of the Clockwork RNN [53] which flattens the network, simply forcing the various hidden states to be updated at different predefined rates. Further approaches in this direction attempt to have the network learn the time scale at which it should update – recent examples include Hierarchical Multiscale RNNs [14] and Gated Feedback RNNs [16] which suggest various ways of connecting, via additional gated connections, the outputs of various layers in a stacked RNN architecture.

These approaches have been successful, although it is not precisely clear how adding extra multiplicative gates could alleviate the vanishing gradient problem. Secondly, many of these approaches add significant extra structures to the network (typically a large LSTM) to force it to behave in a manner which it was already capable of. That significant extra structures need to be added to be able to reliably train it to express such behaviours seems to suggest that a more fundamental change in the model would be wise.

A closely related approach is augmenting RNNs to add an attention mechanism. This is especially popular in translation tasks, in which a common approach is to use two separate RNNs in an *encoder-decoder* arrangement where the first network encodes the sentence into a fixed length vector and the second expands the vector into a new sentence in a different language [12, 59]. Content-based attention is often used to improve these models, by allowing the decoder access to all previous states of the encoder via a weighted sum. Weights are derived from similarities between each encoder state and a vector produced by the decoder at each step [10, 22].

While this achieves excellent performance at extraordinary scale [103], it also requires computing similarities for every step of the output with every step of the input, which will scale quadratically in sentence length. Consequently, it seems like a useful avenue of research to allow the encoder RNN to learn a more useful reduced representation of the sentence to remove the need for the attention mechanism.

A simpler recent architecture of note is the Unitary Evolution RNN [3] which guarantees the eigenvalues of the recurrent weight matrix have a magnitude of one. While this leads to provably non-vanishing or exploding gradients, in practice they still seem to struggle to learn to store information for very long time periods. This is potentially due to the seemingly ad-hoc composition of unitary operators used to allow unconstrained optimisation of parameters while maintaining the desired properties of the recurrent matrix.

Another principled way to approach the design of RNNs is to use *strongly typed* architectures. [4] This approach is inspired by the desire to alter the main source of problems in classic RNNs – the recurrent weight matrix. The argument is that continually applying the same weight matrix to the hidden states leads to incoherent features in the hidden state. The solution is to attempt to *diagonalise* the state updates. This corresponds to ensure all operations on the hidden state apply elementwise, and that the input should be the only item projected through weight matrices.

Another line of enquiry which can help with long time-dependencies focuses on the initialisation of the network. IRNNs [56], which simply initialise the recurrent weights matrix to the identity (presaged by the nearly diagonal initialisation in [67]), perform remarkably well on pathological tasks. The importance of the initialisation of the recurrent weights matrix was further emphasised in [38] where marked differences were found between initialising a

slightly modified vanilla RNN with the identity matrix or a random orthogonal matrix. While simply initialising the network to have certain beneficial properties provides no guarantees on final performance after training, it is important to note the significant results reported especially as they are likely to at least be partially transferable to any given architecture.

Algorithmic Solutions

The second class of attempts to address the issue focuses on techniques for training the network. The first of these is to use approximate second-order methods to try and rescale the gradients appropriately. A number of notable attempts used Hessian-Free Optimization [62, 8] and achieved remarkable results. Perhaps more interestingly, it was subsequently found that many of the results could in fact be achieved with very careful tuning of standard stochastic gradient descent with momentum and careful initialisation [92].

Another interesting approach which seems to help learn networks learn tasks requiring longer time dependencies is to add a regularisation term into the loss to penalise the difference between successive hidden states [54]. This encourages the network to learn smooth transitions and may help it “bootstrap” itself past the vanishing gradients.

A notable recent work applies the recently proposed path-SGD [70] to Vanilla RNNs with ReLU activations. This involves re-casting the RNN as a very deep feed-forward network with shared weights across times-steps. The authors then regularise the update based on the scale of the individual paths through the network. This allows them to achieve significant results with an architecture that is otherwise very challenging to train – the linear rectifier is unbounded in one direction which typically leads rapidly to exploding gradients [71].

2.4.2 Memory

This line of research attempts to enhance the capabilities of RNNs by focusing on the way they interact with their memory. There are two key motivations for this kind of work. One aim is to decouple the memory from the network so that it can store and interact with arbitrary quantities of information. The second is related to vanishing gradients in the hope that a novel method of storing or addressing information will be more robust to long time delays.

A notable effort to decouple RNNs from their memory comes in Neural Turing Machine proposed in [30]. This uses a neural network (which may or may not be recurrent itself) to attend to a large matrix of memory. Addressing for both reading and writing is done via a address mechanism mostly based on the softmax. It was evaluated on a number of synthetic tasks which mostly involved accumulation such as remembering a sequence of arbitrary patterns and outputting them after a number of time steps in a particular order.

Subsequently, Grefenstette et al. propose a set of neural Deques, Queues and Stacks which attempt to address the limitations of the fixed size memory of the Neural Turing Machine [32]. A similar approach is also found in Joulin and Mikolov’s Stack Augmented RNN [47]. These allow for unbounded memory with efficient access, at the cost of the random access nature of the memory. Primarily these are proposed as augmentations to an RNN, with the idea being that a standard RNN could operate as a small working memory while the more complex data structure (with more awkward access semantics) can be used for long-term, more exact storage. This is an exciting direction as it enables something akin to a learnable Von Neumann architecture [30].

Other such approaches have also arisen such as the Neural Random-Access Machine, which is similar to a Neural Turing Machine except that it solves a number of the inefficiencies with a clever addressing mechanism [55]. Other related works include Pointer Networks [100] and Neural GPUs [49]. These tend to involve adding fairly complex structures on

top of an LSTM. While they are highly successful on synthetic tasks when compared to an un-augmented LSTM, very few of the tasks seem like the kind of task that an RNN should inherently be incapable of solving. Nearly all of the tasks addressed require a fixed size memory, and solutions can be imagined composed of carefully writing to it and reading from it. Hence it seems as though it is worth revisiting the basic architecture, rather than searching for distinct modules to add on.

An approach in this direction is Danihelka et al.'s Associative LSTM [19]. This work incorporates an associative memory model into the LSTM by adjusting several of the operations to reflect something close to a Holographic Reduced Representation [79]. This is a very encouraging angle – to re-think the way in which the RNN uses its memory at a basic level seems like the most efficient and sensible way to address the issues. The results were very promising on some interesting synthetic tasks, focusing on the ability of the network to do key-value retrieval which should be a strength of the Holographic Reduced Representation. Unfortunately simply adjusting the network to use such a reduced representation was insufficient for good performance and a number of complicating factors had to be incorporated, including keeping many redundant copies of memory which diminishes somewhat the elegance of the solution.

Many of these ideas seem to be founded on an intuition which has not been clearly formalised. In general, the task of an RNN is to learn a mapping from one sequence to another sequence one step at a time using an intermediate memory (the hidden states). Intuitively a solution must consist of writing to and reading from memory, both conditioned on input. In an RNN, reading from the memory is trivially attended to as typically the entire memory is used to produce the output at each time step. However, we would suggest identifying two separate types of writing to memory: accumulating and forgetting. Many of the architectures discussed below (and the GRU and LSTM above) have separate mechanisms for these two tasks. We consider an accumulation to be an updating of the state, typically additively, to incorporate new information. Forgetting is simply wiping clean a section of the memory – this tends to be achieved multiplicatively. Different tasks will require a different balance of the two. Instinctively, a classification tasks in which the network is only evaluated on its final output would favour more accumulative solutions. Conversely, tasks which may not even have a fixed length but rather require constant prediction should be better suited to networks with an ability to forget information that is no longer relevant.

2.5 Tensors in Neural Networks

The final set of related works provide the key motivation for the next section. In general, many of the above architectures contain multiplicative gates. Multiplicative gating structures are not novel, introduced by Hinton as early as 1981 [40]. Sigaud et al. separate modern uses into two categories based on the intuition behind the application [85]. The first class seeks to use the gate in a manner akin to a transistor in a digital circuit, seeking to switch the flow of information between computational units on or off. This captures the way such connections are used in LSTMs as well as their feed-forward cousins Highway Networks [89]. The second class of gated networks aims to implement a multiplicative connection between inputs. The latter class strictly generalises the former. Intuitively, the former class captures those models where two inputs to a computational unit are multiplied before they are processed, while the second admits some transformations before the multiplication is realised. Two excellent reviews can be found in [63, 85] and rather than repeat in detail the history we attempt to bring out the most salient examples which make clear the intuitions behind incorporating multiplicative interactions. Finally, we mention a few particularly very recent applications

that have achieved significant successes after publication of the aforementioned reviews.

The earliest examples of multiplicative interactions in the context of the current Deep Learning movement is in the Gated Boltzmann machine [64] which generalised Restricted Boltzmann Machines [87] (a class of probabilistic graphical models very closely related to neural networks) to model conditional relationships between pairs of images. These parameterise a multiplicative connection with a three-way tensor \mathbf{W} . The likelihood of a hidden state being 1 given a pair of inputs \mathbf{x} and \mathbf{y} is defined as

$$p(h_k|\mathbf{x}, \mathbf{y}) = \sigma \left(\sum_{ij} W_{ijk} x_i y_j \right) \quad (2.5)$$

where $\sigma(\cdot)$ is the sigmoid discussed earlier [64]. This architecture very naturally extends the operation of a feed-forward network to deal with two inputs (compare equations (2.1) and (2.5)), but the size of the tensor \mathbf{W} is a serious limitation. This work outlines three key points: multiplicative interactions are somehow a natural way of handling multiple distinct inputs, multiplicative interactions are properly parameterised by tensors and tensors can be prohibitively large. This last point is addressed in later works in the same field by factorising the tensor into three matrices [95, 65].

The Multiplicative RNN [91, 90] picks up both the notion that multiplicative interactions are sensible ways to model interactions between two inputs and the factorisation from the above. Applied in a language modelling context, the reasoning is that it makes sense to allow the current input symbol to affect the recurrent transition matrix. These MRNNs achieved significant results, especially when combined with Hessian Free [62], a method of approximating second order optimisation.

Very recent methods have revisited a number of these ideas, although often with little to no mention of the underlying tensorial structure. Of note is the Multiplicative Integration RNN which proposes adjusting the building block of eq. (2.2) to use element-wise multiplication instead of addition [104]. This corresponds to factorising the tensor into two matrices and has an extreme gating effect, especially during back-propagation. In order to successfully learn, it was found necessary to add additional bias terms to the point where the final model required four weight matrices and five bias vectors per unit. In some experiments they incorporate this approach into an LSTM, leading to a model which would have nearly 16 times the parameters of the simplest RNN proposed by eq. (2.2). We suggest that paying more attention to the tensorial structure of the desired interaction would lead to more thoroughly grounded models, and propose one such model in chapter 4.

LSTM-style gating also continues to receive attention, being employed in Gated-Attention Readers to combine representations of questions and potential answers [20]. It is also used in PixelCNN and WaveNet [74, 73], allowing to model to choose how particular features propagate upwards through layers.

It is important to realise that all of these multiplicative structures have an inherent tensor structure, whether it is made implicit or not. In chapter 3 we investigate with more rigour the types of tensor products which generalise many of these multiplicative interactions.

Chapter 3

Three-way Tensors and Bilinear Products

In order to represent functions of two arguments with neural networks, three-way tensors are unavoidable. In this chapter we discuss the vector-tensor-vector bilinear product and assess its suitability, both theoretically and empirically, to be used inside a neural network,

This product admits a surprising number of interpretations including theorem 1 – that it allows us to operate on a pairwise exclusive-or of features. This is a powerful result which shows that admitting tensors into neural networks can solve the exclusive-or problem without any hidden units. Further, we demonstrate the expressive power of the tensor by showing that a polynomially sized network with a single tensor product is capable of expressing a class of functions which standard two layer networks can not approximate accurately without exponential size.

The downside of these tensor products is the storage requirements. We investigate methods for decomposing tensors into products of smaller objects. Finally we conduct experiments to show that learning the coefficients of a decomposed tensor can be done in situ with gradient descent. We attempt to investigate the tradeoff between the complexity of the decomposition and the expressivity of the product and find that it scales as expected.

3.1 Definitions

Formally, we refer to a multi-dimensional array which requires n indices to address a single (scalar) element as a n -way tensor and occasionally refer to n as the number of dimensions of the tensor. In this sense a matrix is a two-way tensor and a vector a one-way tensor, although we will use their usual names for clarity. Notationally, we attempt to stick to the notation used in [52] deviating only where it would become unwieldy. We denote tensors with three or more dimensions with single calligraphic boldface letters such as \mathcal{W} . Matrices and vectors will be denoted with upper and lower case boldface letters while scalars will be standard lower case letters. We will often need to address particular substructures of tensors. This is analogous to pulling out individual rows or columns of a matrix. To perform this we fix some of the indices to specific values and allow the remaining indices to vary across their range. We denote by \cdot the indices allowed to vary, the rest will be provided with a specific value. For example, $\mathbf{A}_i \cdot$ would denote the i -th row of the matrix \mathbf{A} . For three-way tensors we refer to the vector elements produced by fixing two indices as *threads*. It is also possible to form matrices by fixing only one index – we refer to these as *slices*. Table 3.1 provides examples of the possibilities.

When dealing with tensor-tensor products, it is important to be precise as there are often

<i>Description</i>	<i>Example</i>
scalar	a
vector	\mathbf{b}
matrix	\mathbf{C}
higher order tensor	\mathcal{D}
element of vector (scalar)	b_i
element of matrix (scalar)	C_{ij}
element of 3-tensor (scalar)	D_{ijk}
row of matrix (vector)	$\mathbf{C}_{i\cdot}$
column of matrix (vector)	$\mathbf{C}_{\cdot i}$
<i>fiber</i> of 3-tensor (vector)	$\mathcal{D}_{\cdot jk}$
<i>slice</i> of 3-tensor (matrix)	$\mathcal{D}_{\cdot\cdot k}$

Table 3.1: Example of notation for tensors.

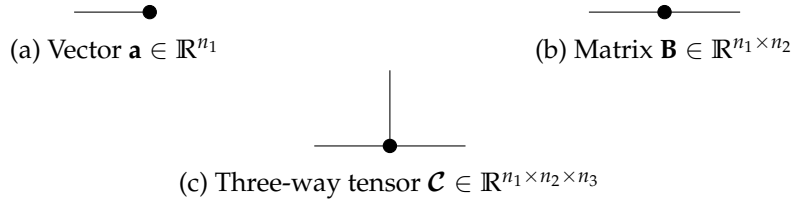


Figure 3.1: Example tensor network diagrams.

a number of possible permutations of indices that would lead to a valid operation. The downside of this is that it leads to unwieldy notation. Fortunately, we are only concerned with a couple of special cases. In particular, we need to multiply a three-tensors by vectors and a matrices by vectors. Matrix-vector multiplication consists of taking the dot product of the vector with each row of the matrix. For example, with a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a vector $\mathbf{x} \in \mathbb{R}^n$, if $\mathbf{y} = \mathbf{A}\mathbf{x}$ (with \mathbf{y} necessarily in \mathbb{R}^m), then

$$y_i = \sum_j^n A_{ij} x_j = \langle \mathbf{A}_{i\cdot}, \mathbf{x} \rangle$$

where $\langle \cdot, \cdot \rangle$ denotes the inner (dot) product. This can be viewed as taking all of the vectors formed by fixing the first index of \mathbf{A} while allowing the second to vary and computing their inner product with \mathbf{x} . To perform the same operation using the columns of \mathbf{A} we need to fix the second index, the would typically be done by exchanging the order: $\mathbf{x}^\top \mathbf{A}$. This kind of operation is sometimes referred to as a *contractive* product, especially in the physical sciences [75]. This name arises because we form the output by joining a shared dimension (by element-wise multiplication) and contracting it with a sum.

We can generalise the operation to tensors: choose an index over which to perform the product, collect every thread formed by fixing all but that one index and compute their inner product with the given vector. If the tensor has n indices, the result will have $n - 1$. A three tensor is in this way reduced to a matrix. Kolda and Bader introduce the operator $\bar{\times}_i \cdot$ for this, where i represents the index to vary [52]. For the bilinear forms we are concerned with, this leads to the following notation:

$$\mathbf{z} = \mathcal{W} \bar{\times}_1 \mathbf{x} \bar{\times}_2 \mathbf{y}. \quad (3.1)$$

When \mathcal{W} is a three-way tensor, we prefer a more compact notation

$$\mathbf{z} = \mathbf{x}^\top \mathcal{W} \mathbf{y}. \quad (3.2)$$

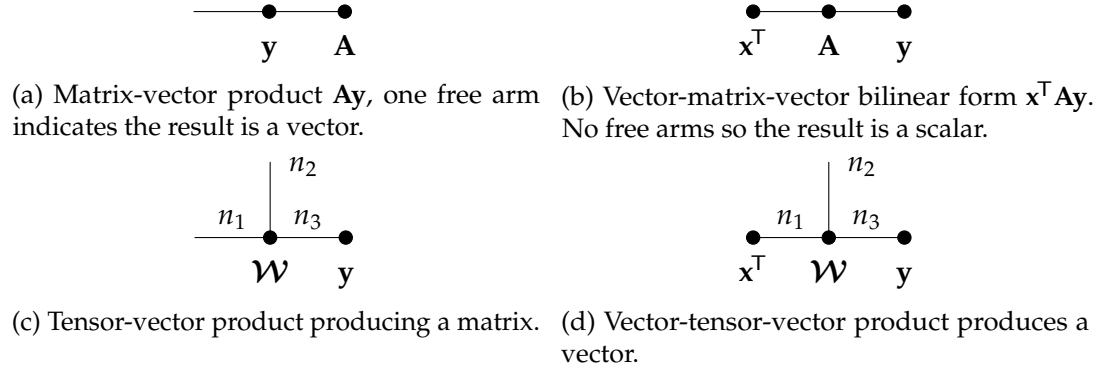


Figure 3.2: Various products expressed as Tensor Network Diagrams.

This loses none of the precision of the more verbose notation provided we make clear that we intend \mathbf{x} to operate along the first dimension of the tensor and \mathbf{y} the third such that $(\mathbf{x}^T \mathcal{W})\mathbf{y}$ exactly corresponds with equation (3.1). With either notation, for a tensor $\mathcal{W} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, we must have that $\mathbf{x} \in \mathbb{R}^{n_1}$, $\mathbf{y} \in \mathbb{R}^{n_3}$ and the result $\mathbf{z} \in \mathbb{R}^{n_2}$.

An intuitive way to illustrate these ideas is using Tensor Network Diagrams [17, 75]. In these diagrams, each object is represented as a circle, with each free ‘arm’ representing an index used to address elements. A vector therefore has one free arm, a matrix two and so on. Scalars will have no arms. Figure 3.1 provides examples for these simple objects.

Where these diagrams are especially useful is for representing contractive products where we sum over the range of a shared index. We represent this by joining the respective arms. As an example, a matrix-vector product $\mathbf{y} = \mathbf{A}\mathbf{x}$ has such a contraction: $y_i = \sum_j A_{ij}x_j$. This is shown in figure 3.2a – it is clear that there is only a single free arm, so the result is a vector as it should be. Figure 3.2 shows some examples of these kinds of products.

3.2 Interpreting Bilinear Products

There are several ways to describe the operation performed by the bilinear products we are concerned with. These correspond to different interpretations of the results. The following interpretations provide insight both into what is actually being calculated and how the product might be applicable to a neural network setting. In all of the below we use the above definitions of \mathbf{x} , \mathbf{y} , \mathbf{z} and \mathcal{W} .

3.2.1 Stacked bilinear forms

If we consider the expression for a single element of \mathbf{z} , we get

$$z_j = \sum_i^{n_1} \sum_k^{n_3} W_{ijk} x_i y_k$$

as expected. We can re-write this in terms of the slices of \mathcal{W} :

$$z_j = \mathbf{x}^T \mathcal{W}_{\cdot j} \mathbf{y}$$

which reveals the motivation behind the notation in equation (3.2). It also reveals that each element of \mathbf{z} is itself linear in \mathbf{x} or \mathbf{y} if the other is held constant.

This provides an interpretation in terms of similarities. If we consider the standard dot

product of two vectors \mathbf{a} and \mathbf{b} of size m :

$$\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{b} = \sum_{i=1}^m a_i b_i = \cos \theta \cdot \|\mathbf{a}\|_2 \cdot \|\mathbf{b}\|_2$$

where θ is the angle in the angle between the vectors. If the product is positive, the two vectors are pointing in a similar direction and if it is negative they are in opposite directions. If it is exactly zero, they must be orthogonal. The dot product therefore provides us with a notion of the similarity between two vectors. Indeed if we normalise the vectors by dividing each component by their l_2 norm we recover exactly the widely used cosine similarity, common in information retrieval [86, 94]. Note that we can generalise this idea by inserting a matrix of (potentially learned) weights \mathbf{U} which enables us to define general scalar bilinear forms

$$\langle \mathbf{a}, \mathbf{U}\mathbf{b} \rangle = \langle \mathbf{a}^T \mathbf{U}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{U} \mathbf{b}.$$

In a bilinear tensor product, each component of the result takes this form. We can therefore think of the product as computing a series of distinct similarity measures between the two input vectors. With this in mind the obvious question is: what is the role of the matrix? The first thing to note is that inserting a matrix into the inner product allows the two vectors to be of different dimension. We also observe that a matrix-vector multiplication consists of taking the dot product of the vector with each row or column of the matrix. Given our current interpretation of the dot product as an un-normalised similarity measure, we can also interpret a matrix-vector multiplication as computing the similarity of the vector with each row or column of the matrix.

We can then think of the rows of the matrix \mathbf{U} in the above as containing patterns to look for in the \mathbf{b} vector and the columns to contain patterns to test for in the \mathbf{a} vector. If we consider the vectors \mathbf{a} and \mathbf{b} to come from different feature spaces, the matrix \mathbf{U} provides a conversion between them allowing us to directly compare the two. We can then interpret each coordinate of the result of the bilinear tensor product as being an independent similarity measure based on different interpretations of the underlying feature space. In this sense, where a matrix multiplication looks for patterns in a single input space, a bilinear product looks for *joint* patterns in the combined input spaces of \mathbf{x} and \mathbf{y} . This becomes clear if we write out each element of the result vector

$$\mathbf{z} = \begin{bmatrix} \mathbf{x}^T \mathbf{W}_{\cdot 1} \mathbf{y} \\ \mathbf{x}^T \mathbf{W}_{\cdot 2} \mathbf{y} \\ \vdots \\ \mathbf{x}^T \mathbf{W}_{\cdot n_2} \mathbf{y} \end{bmatrix} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{W}_{\cdot 1} \mathbf{y} \rangle \\ \langle \mathbf{x}, \mathbf{W}_{\cdot 2} \mathbf{y} \rangle \\ \vdots \\ \langle \mathbf{x}, \mathbf{W}_{\cdot n_2} \mathbf{y} \rangle \end{bmatrix} = \begin{bmatrix} \cos \theta_1 \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{W}_{\cdot 1} \mathbf{y}\|_2 \\ \cos \theta_2 \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{W}_{\cdot 2} \mathbf{y}\|_2 \\ \vdots \\ \cos \theta_{n_2} \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{W}_{\cdot n_2} \mathbf{y}\|_2 \end{bmatrix}.$$

3.2.2 Choosing a matrix

Following on from the above discussion we claim that for each coordinate of the output we are computing a *similarity vector* which we compare to the remaining input to generate a scalar value. If we consider all coordinates at once, we see that this amounts to having one input choose a matrix, which we then multiply by the remaining input vector. We have refrained from making these points in terms of the specific vectors referenced above to make the point that the operation is completely symmetrical. While it aids interpretation to think of one vector choosing a matrix for the other vector, we can always achieve the same intuition after switching the vectors, give or take some transposes.

This interpretation is very clear from the expression of the product in equation (3.1). Simply by inserting parentheses we observe that we are first generating a matrix in a way

somehow dependent on the first input and multiplying the second input by that matrix. In this sense we allow one input to choose patterns to look for in the other input.

This intuition of choosing a matrix is suggested in [90] in the context language modelling with RNNs. It is suggested that allowing the current input character to choose the hidden-to-hidden weights matrix should confer benefits. This intuition (and the factorisation of the implicit tensor) was put to use earlier in of Conditional Restricted Boltzmann Machines [95] which actively seek to model the conditional dependencies between two types of input.

Although this provides a powerful insight into the bilinear product, it is worth reinforcing that the product is entirely symmetrical. We can not think purely about it as x choosing a matrix for y as the converse is equally true.

3.2.3 Tensor as an independent basis

Extending the above to try and capture the symmetry of the operation, we introduce the notion that the coefficients of the tensor represent a basis in which to compare the two inputs, independent of both of them. This idea is mentioned in [97] which considered the problem of data that can be described by two independent factors. The tensor then contains a basis which characterises the interaction by which the factors in the input vectors combine to create an output observation.

This is a somewhat abstract interpretation – to attempt to create a more concrete intuition consider the case when both vectors have a single element set to 1 and the remainder 0. In this case the first tensor-vector product corresponds to taking a slice of the tensor, resulting in a matrix. The final matrix-vector product corresponds to picking a row or column of the matrix. Consequently the whole operation is precisely looking up a fibre of the tensor. To generalise from such a “one-hot” encoding of the inputs to vectors of real coefficients, we simply replace the idea of a *lookup* with that of a *linear combination*. The vectors then represent the coefficients of weighted sums; first over slices and then over rows or columns. The final product is then a distinct representation of both vectors in terms of the independent basis expressed by the tensor.

3.2.4 Operation on the outer product

Perhaps the most powerful interpretation of the product is achieved by observing that it is a linear operation on pairwise products of inputs. This interpretation is essential for understanding the expressive power of the operation as it gives rise to an obvious XOR-like behaviour. It also serves as a useful reminder that a bilinear form is linear in each input only when the other is held constant – when both are allowed to vary we can represent complex non-linear relations. A way to approach this interpretation arises from a method of implementing the bilinear product in terms of straightforward matrix operations.

To discuss this we need to introduce the *matricisation* of a tensor. Intuitively this operation is a way of *unfolding* or *flattening* a tensor into a matrix, preserving its elements. Specifically the mode- n matricisation of a tensor is defined as an operation which takes all mode- n fibres of the tensor and places them as columns to create a matrix. We denote an mode- n matricisation of a tensor \mathcal{W} as $\text{mat}_n(\mathcal{W})$. While the operation is fairly straightforward, describing the exact permutation of the indices is awkward – for a robust treatment of the general case (and the source of the above definition) see [52].

As an example consider the three-way tensor $\mathcal{X} \in \mathbb{R}^{2 \times 3 \times 4}$ with slices

$$\begin{aligned} \mathbf{X}_{1..} &= \begin{bmatrix} a & d & g & j \\ b & e & h & k \\ c & f & i & l \end{bmatrix} \\ \mathbf{X}_{2..} &= \begin{bmatrix} m & p & s & v \\ n & q & t & w \\ o & r & u & x \end{bmatrix}. \end{aligned} \quad (3.3)$$

To construct $\text{mat}_2(\mathcal{X})$ we fix the first and third indices and place the resulting threads as columns of a new matrix. Fixing the first index corresponds to choosing one of the slices presented above while fixing the third amounts to choosing a column of one of the slices. The result is:

$$\text{mat}_2(\mathcal{X}) = \begin{bmatrix} a & d & g & j & m & p & s & v \\ b & e & h & k & n & q & t & w \\ c & f & i & l & o & r & u & x \end{bmatrix}. \quad (3.4)$$

Although this notion captures and generalises the vectorisation operator often encountered in linear algebra, we retain the classical vec operator for clarity. This flattens a matrix into a vector by stacking its columns. For some matrix \mathbf{A} with n columns:

$$\text{vec}(\mathbf{A}) = \begin{bmatrix} \mathbf{A}_{\cdot 1} \\ \mathbf{A}_{\cdot 2} \\ \vdots \\ \mathbf{A}_{\cdot n} \end{bmatrix}.$$

For our purposes it is sufficient to note that a mode-2 matricisation of a three-way tensor $\mathcal{W} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ must have shape $n_2 \times n_1 n_3$. The following lemma helps us describe the components of this matricisation.

Lemma 3.2.1 (Matricisation/vectorisation). *The j -th row of the mode-2 matricisation of the three-way tensor \mathcal{W} is equivalent to the vectorisation of the transpose of the slice formed by fixing the second index at j :*

$$\text{mat}_2(\mathcal{W})_{j\cdot} = \text{vec}(\mathbf{W}_{\cdot j}^\top).$$

Proof. By the above definition of the vectorisation operator, each index (i, k) in some matrix $\mathbf{U} \in \mathbb{R}^{n_1 \times n_3}$ maps to element $i + (k - 1)n_1$ in $\text{vec}(\mathbf{U})$. By the definition of the mode-2 matricisation we would expect to find tensor element W_{ijk} at index $(j, i + (k - 1)n_3)$. Hence if we fix j , we have precisely the vectorisation of the j -th slice of the tensor.

The indices into $\text{mat}_2(\mathcal{W})$ can be thought of as arising from the following construction. First fix all indices to 1. Construct a column by sweeping the second index, j , through its full range. Then increment the first index i and repeat the procedure, placing the generated columns with index i . Only when i has swept through its full range increment the final index k and repeat the procedure. The generated columns should then be at positions $i + (k - 1)n_1$. \square

To understand the lemma it helps to consider an example. Using the $2 \times 3 \times 4$ tensor \mathcal{X} defined in equation (3.3), fixing the second index to 2 gives a 4×2 transposed slice:

$$\mathbf{X}_{2\cdot}^\top = \begin{bmatrix} b & n \\ e & q \\ h & t \\ k & w \end{bmatrix}.$$

Vectorising this matrix takes each column and stacks them, giving

$$\text{vec}(\mathbf{X}_{\cdot 2}^T) = \begin{bmatrix} b \\ e \\ h \\ k \\ n \\ q \\ t \\ w \end{bmatrix}$$

which is precisely row two in equation (3.4).

These flattenings are important as they allow us to implement many operations involving tensors in terms of a small number of larger matrix operations when compared to the naive approach.

Lemma 3.2.2 (Matricised product). *For a tensor \mathcal{W} and vectors \mathbf{x}, \mathbf{y} as above, we can describe the product $\mathbf{z} = \mathbf{x}^T \mathcal{W} \mathbf{y}$ in terms of the mode-2 matricisation of \mathcal{W} as follows:*

$$\mathbf{z} = \text{mat}_2(\mathcal{W}) \text{vec}[\mathbf{y} \mathbf{x}^T] \quad (3.5)$$

Proof. To prove this we can compare the expressions for a single element of the result. An element z_j from equation (3.5) is formed as the inner product of the j -th row of the flattened tensor and the vectorised outer product of the inputs. By lemma 3.2.1:

$$z_j = \sum_{s=1}^{n_1 n_3} (\text{mat}(\mathcal{W}))_{js} \cdot \left(\text{vec}[\mathbf{y} \mathbf{x}^T] \right)_s$$

We replace the sum over s with a sum over two indices, i and k , and using them to appropriately re-index the flattenings as described in lemma 3.2.1 we derive

$$\begin{aligned} z_j &= \sum_{i=1}^{n_1} \sum_{k=1}^{n_3} W_{ijk} x_i y_k \\ &= \mathbf{x}^T \mathbf{W}_{\cdot j} \mathbf{y} \end{aligned}$$

□

Therefore to understand the bilinear product it helps to understand the matrix $\mathbf{y} \mathbf{x}^T$. This matrix takes the following form:

$$\mathbf{y} \mathbf{x}^T = \begin{bmatrix} x_1 y_1 & \dots & x_{n_1} y_1 \\ \vdots & \ddots & \vdots \\ x_1 y_{n_3} & \dots & x_{n_1} y_{n_3} \end{bmatrix}$$

which contains all possible products of pairs of elements, one from each vector. This captures some interesting interactions.

Theorem 1 (Bilinear exclusive-or). *Bilinear tensor products operate on the pair-wise exclusive-or of the signs of the inputs.*

Proof. Consider the sign of a scalar product $c = a \cdot b$. If one of the operands a or b is positive and the other negative, then the sign of c is negative. If both are positive or both are negative, then the result is positive. This captures precisely the “one or the other but not both” structure of the exclusive-or operation.

By lemma 3.2.2 a bilinear product can be viewed as a matrix operation on the flattened outer product of the two inputs. As each element in the outer product is the product of two scalars, the signs have an exclusive-or structure. \square

Corollary 1.1 (Bilinear conjunctions). *If the inputs are binary, bilinear products operate on pairwise conjunctions of the inputs.*

Proof. If $a, b \in \{0, 1\}$, then $a \cdot b = 1$ if and only if both a and b are 1. If either or both are 0 then their product must be zero. Following the same structure as the proof of theorem 1, for binary inputs we must have this conjunctive relationship. \square

Implications

Firstly, this captures the notion that the bilinear tensor product operates implicitly on higher level features, constructed by combining both inputs. This indicates that it is capable of capturing complex binary relationships based on correlations between both inputs. This allows it to model a rich class of binary relations between \mathbf{x} and \mathbf{y} .

Secondly this suggests considering the case where both inputs are the same: $\mathbf{z} = \mathbf{x}^\top \mathcal{W} \mathbf{x}$. Replacing bilinear forms with quadratic forms may invalidate some of the similarity based interpretations, but it also provides an interesting building block for feed-forward networks. As an example, note that a plain perceptron has been long known to be incapable of learning the exclusive-or mapping [68] – a neural network to solve the problem requires either a hidden layer or an additional input feature (specifically the conjunction of the inputs) [82]. By corollary 1.1 this quadratic tensor form would implicitly and naturally capture the additional conjunctive feature and be capable of solving the exclusive or problem without hidden layers or hand-engineered features.

Finally we will point out that this gives an indication of what we term the ‘apparent depth’ of the tensor product. While it has long been known that neural networks with a single hidden layer and appropriate non-linearity (and therefore two weight matrices) are universal function approximators [45], recent results have suggested that the depth of the network is important in keeping the number of parameters bounded [23, 96]. In particular, Eldan and Shamir show that for a particular class of radial basis-like functions – which depend only on the squared norm of the input ¹ – there exist networks with two hidden layers and a number of nodes polynomial in the input dimension which can perfectly approximate the function. They then show that a shallower network with a single hidden layer can only perfectly approximate the function with a number of nodes exponential in the input dimension [23]. This is a very recent result (published in COLT 2016), which is remarkable as it proves an impossibility: there is a realistic class of functions which it is impossible to approximate accurately with a two layer network with polynomial size. However, this restriction can be subverted if we allow tensor semantics in the first layer of the network.

To show this, we first require the following lemma.

Lemma 3.2.3 (Squared norm). *A tensor layer $\mathbf{x}^\top \mathcal{W} \mathbf{x}$ is capable of exactly computing the squared Euclidean norm $\|\mathbf{x}\|_2^2 = \sum_i x_i^2$.*

Proof. Recalling lemma 3.2.2

$$\mathbf{x}^\top \mathcal{W} \mathbf{x} = \text{mat}_2(\mathcal{W}) \text{vec}(\mathbf{x} \mathbf{x}^\top).$$

¹Precisely, functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $f(\mathbf{x}) = f(\mathbf{y})$ implies that $\|\mathbf{x}\|^2 = \|\mathbf{y}\|^2$. Eldan and Shamir deal specifically with cases where the norm is the Euclidean norm (as is the case for the common squared-exponential kernel). They also require some additional constraints on the boundedness and support of the function, see [23] theorem 1, for the precise construction.

Note that the diagonal elements of $\mathbf{x}\mathbf{x}^\top$ are the components x_i^2 which we need to sum to compute the squared norm. We simply need to define $\text{mat}_2(\mathcal{W})$ such that it picks out the correct elements and sums them. If $\mathbf{x} \in \mathbb{R}^d$ then $\mathcal{W} \in \mathbb{R}^{d \times 1 \times d}$ as we only require a single output. The matricisation is then a $1 \times d^2$ row vector and we simply need to set its elements zero apart from those with indices of the form $1 + id$ for $i \in \{1, \dots, d\}$. The matrix product will then simply pick out the diagonal elements and sum them, computing the squared norm as required. \square

Theorem 2 (Exponential expressivity of tensor layers). *There exists a class of functions which can be approximated arbitrarily well by a two layer network including a single tensor layer with a number of nodes polynomial in the size of the input which require an exponential number of nodes to be approximated by a standard feed-forward layer.*

Proof. We make use of Eldan and Shamir’s proof [23]. This is highly technical so we only provide a constructive outline here, which shows where we need to fit in. The proof in [23] shows that for a given class of radial basis-like functions a two-layer feed-forward network requires exponential width to be able to approximate below a certain error.

They then construct a three-layer network in the following manner: for an input \mathbf{x} , the first two layers learn to approximate the map $\mathbf{x} \mapsto \|\mathbf{x}\|_2^2 = \sum_i x_i^2$ while the final layer learns a univariate function on the squared norm, it is then proved that this required that each layer have polynomial width.

Using lemma 3.2.3, we can compute the squared norm exactly with a single layer. We can therefore construct a network in the same fashion as Eldan and Shamir, replacing the first two layers with a single tensor layer.

This plugs directly into the proof of lemma 10 in [23, pp. 18–19], simply noting that we must duplicate the tensor layer a polynomial number of times. At most we will need the same number of nodes in the final hidden layer of Eldan and Shamir’s three layer network. \square

This is a remarkable result, as it shows the increase in expressive power gained by allowing tensor representations. We have shown a single tensor layer can do the work of two feed-forward layers. This provides a strong motivation to explore architectures incorporating these tensor semantics.

3.3 Tensor Decompositions

Bilinear products show great theoretical promise. Unfortunately such a tensor is often very large; an $n \times n \times n$ would require $\Theta(n^3)$ space to store explicitly and a cubic number of operations to use to compute a bilinear product. Further, it is highly likely that we will be looking to model simpler interactions than the full tensor might represent. This leads to the idea of a parameterised decomposition – an ideal solution would be a method of representing the tensor with some way of specifying the complexity of the bilinear relationship and hence making explicit the relationship between the number of parameters required and the expressive power of the model.

We will investigate two methods of tensor decomposition with a view towards their use as a building block for neural networks. The outcomes we wish to evaluate are: the savings in storage, how convenient it is to specify the number of parameters, convenience and efficiency of implementation (in terms of matrix operations) and whether the gradients suggest any benefits or hindrances when learning by gradient descent. There is a significant amount of research into tensor decompositions, [52] provides a thorough review. Here we consider some of the most general decompositions and families of decompositions without going into detail outside of where it is relevant to the above concerns.

3.3.1 CANDECOMP/PARAFAC

This method of decomposing a tensor dates as far back as 1927 [43, 42] and was rediscovered repeatedly across a range of communities over the next 50 years. [52] First referred to as ‘polyadic form of a tensor’ we refer to it as the CP-decomposition after two prominent publications in 1970 referring to the technique as ‘canonical decomposition’ (CANDECOMP) [9] or ‘parallel factors’ (PARAFAC). [35] It is a fairly straightforward extension of a matrix rank decomposition to the more general case, representing a tensor as a sum of rank one tensors. This section contains a formal description of the three-way case as well as a note on computing bilinear products with the decomposed tensor. Some interesting results on the uniqueness of the decompositions can be found in [52], but they are not relevant to the discussion here.

Description

Given a three-way tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ we wish to approximate it as

$$\mathcal{X} \approx \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$

where R is the rank of the decomposition and $\cdot \otimes \cdot$ denotes the tensor product. This product expands the dimensions; for the first two vectors it is the outer product $\mathbf{a}\mathbf{b}^T$ which generates a matrix consisting of the product of each pair of elements in \mathbf{a} and \mathbf{b} . With three such products we proceed analogously, except now our structure must contain the product of each possible triple of elements. Hence we need three indices to address each entry, so it best described as a three-way tensor. Each individual element has the form

$$X_{ijk} = \sum_{r=1}^R a_{ri} b_{rj} c_{rk}.$$

It is convenient to stack these factors into matrices – we differ slightly from [52] by defining the *factor matrices* of the form $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_R]^T \in \mathbb{R}^{R \times n_1}$ (and equivalently for \mathbf{B} and \mathbf{C}) such that the *rows* of the matrix correspond to the R vectors. We denote a tensor decomposed in this format $\mathcal{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{CP}$.

Bilinear Product

We wish to compute a product of the form $\mathbf{z} = \mathbf{x}^T \mathcal{W} \mathbf{y}$ where $\mathcal{W} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{CP}$ is represented as a CP-decomposition. Conveniently, this can be done with simple matrix products.

Proposition 3.3.1.

$$\begin{aligned} \mathbf{z} &= \mathbf{x}^T \mathcal{W} \mathbf{y} \\ &= \mathbf{B}^T (\mathbf{A} \mathbf{x} \odot \mathbf{C} \mathbf{y}) \end{aligned}$$

when $\mathcal{W} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{CP}$ and \odot represents the Hadamard (element-wise) product.

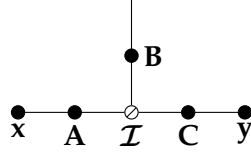


Figure 3.3: One way of representing a bilinear product in the CP-decomposition.

Proof. This follows from straightforward rearranging. Firstly

$$\begin{aligned}
 z_j &= \sum_i^{n_1} \sum_k^{n_3} W_{ijk} x_i y_k \\
 &= \sum_i^{n_1} \sum_k^{n_3} \sum_r^R A_{ri} B_{rj} C_{rk} x_i y_k \\
 &= \sum_r^R B_{rj} \left(\sum_i^{n_1} A_{ri} x_i \cdot \sum_k^{n_3} C_{rk} y_k \right). \tag{3.6}
 \end{aligned}$$

This implies we can compute the quantity inside the brackets for all r at once simply by $\mathbf{Ax} \odot \mathbf{Cy}$ which results in a length R vector. Equation (3.6) then notes that for each output j we take the j -th column of the $R \times n_2$ factor matrix \mathbf{B} and perform a dot product with our earlier result. A series of dot products can be easily expressed with matrix multiplication, noting that we have to transpose \mathbf{B} to keep it on the left but still sum over each column. Hence:

$$\mathbf{z} = \mathbf{B}^T (\mathbf{Ax} \odot \mathbf{Cy}).$$

□

This proposition casts the CP decomposition into a form that looks very similar to some of the RNNs discussed in chapter 2. A key similarity is with the GRU – the manner in which the reset gate is applied to the hidden states during the production of the candidate update \mathbf{z}_t could almost be viewed as a non-linear version of the above. In the LSTM as well there are a number of occasions where very similar calculations are performed.

We can express the bilinear product above as a tensor network diagram, although we have to insert an auxiliary $R \times R \times R$ tensor, denoted \mathcal{I}_R which is defined as

$$I_{ijk} = \begin{cases} 1 & \text{if } i = j = k \\ 0 & \text{otherwise} \end{cases}$$

in a manner analogous to the identity matrix. Taking a bilinear product with this tensor simply represents element-wise multiplication of the two vectors. In the diagram this is represented with the symbol \odot to emphasise the diagonality (and to emphasise that it is different as we do not need to store its parameters). This is presented in figure 3.3, which can be compared to figure 3.2d.

Discussion

The CP decomposition does a remarkable job of preserving the expressivity. To illustrate this consider a special case of the form

$$\mathcal{W} = \sum_r^R \mathbf{a}_r \otimes \mathbf{e}_r \otimes \mathbf{c}_r$$

where the \mathbf{e}_i are the standard basis vectors, with all elements zero apart from position i which is one. Consider $\mathbf{z} = \mathbf{x}^\top \mathcal{W} \mathbf{y}$. Suppose, without loss of generality, that $\mathbf{z} \in \mathbb{R}^R$. Then

$$z_r = \langle \mathbf{a}_r, \mathbf{x} \rangle \cdot \langle \mathbf{b}_r, \mathbf{y} \rangle.$$

If we apply a sigmoid non-linearity, then $\sigma(z_r) > 0.5$ implies z_r is positive so the dot products must have the same sign. $\sigma(z_r) < 0.5$ implies z_r is negative. If inserted in a neural network, the CP decomposition has the ability to compute the exclusive-or of features.

Bilinear products with the CP decomposition also retain the ability to express element-wise multiplication. The tensor it represents must be *diagonal* with the diagonal elements having value one (see proposition A.0.1 in the appendix for a proof). This means all entries ijk such that $i = j = k$ must be one, all other entries zero. We denote such a tensor by $\mathcal{H} \in \mathbb{R}^{N \times N \times N}_{\text{element-wise}}$.

This is straightforward to model with a CP decomposition: $\mathcal{H} = [\mathbf{I}, \mathbf{I}, \mathbf{I}]_{CP}$ where \mathbf{I} is the $N \times N$ identity matrix. This can be easily verified as the definition of a bilinear product with a CP-decomposed tensor includes an element-wise product so it suffices to ensure the factor matrices do not alter the inputs. The CP-decomposition reduces the required parameters to $3N^2$ as opposed to the original N^3 . This is a remarkable reduction, considering that an analogous diagonal *matrix* is the worst case for the corresponding two-way decomposition.

The TT-decomposition is not capable of reducing the parameters. In order to ensure that the central tensor product passes through the appropriate elements it must in fact be equal to \mathcal{H} itself.

3.3.2 Tensor Train

The *tensor-train* decomposition has a much shorter history – it was first proposed in 2011 as a simpler way of representing a slightly hierarchical form derived from a generalisation of the singular value decomposition. [76] It is proposed as an alternative to the CP-decomposition which purportedly provides benefits in terms of numerical stability. It has been used to learn compressed weight matrices in large neural networks [72] suggesting it is a prime candidate to learn with gradient descent.

It is also important to note that the tensor-train can be used more creatively. Novikov et al. [72] use the decomposition to compress the final layers of a large convolutional neural network by reshaping the matrices into high dimensional tensors. The reduction in parameters for the tensor-train is most notable with particularly high-dimensional tensors, so this approach allows for significant compression (compressing one layer up to 200000 times). They also show methods of computing the required matrix vector products without having to expand the tensor. While this approach is highly successful, the implementation (especially of back-propagation) is somewhat involved and it admits a very large number of ways of applying the decomposition. To keep our search space of decompositions feasible, we consider only the straight-forward application of the tensor-train.

Description

The tensor-train decomposition (TT-decomposition) of a general tensor \mathcal{X} with d indices is the tensor $\mathcal{Y} \approx \mathcal{X}$ with elements expressed as a product of slices of three-way tensors (so a product of matrices):

$$Y_{i_1 i_2 \dots i_d} = \mathbf{G}[1]_{\cdot i_1} \cdot \mathbf{G}[2]_{\cdot i_2} \cdots \mathbf{G}[d]_{\cdot i_d} \cdot$$

If dimension j is of size n_j , then $\mathcal{G}[j]$ is size $r_{j-1} \times n_j \times r_j$ so that each slice $\mathbf{G}[j]_{\cdot i}$ is size $r_{j-1} \times r_j$. The collection of r_i is the ‘tt-rank’ of the decomposition and controls the number of

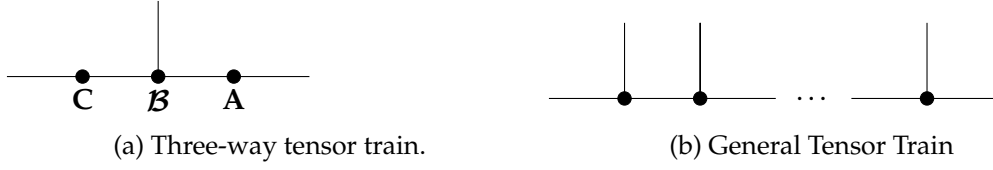


Figure 3.4: Diagrams of the Tensor Train decomposition.

parameters. In order to ensure the result of the chain of matrix products is a scalar, it must be that $r_0 = r_d = 1$. [76]

The three-way case presents no obvious simplifications from the general case but we present it here to ensure consistent notation through the remainder of this section. The TT-decomposition of a three-way tensor of size $n_1 \times n_2 \times n_3$ has three ‘cores’ with shapes $1 \times n_1 \times r_1$, $r_1 \times n_2 \times r_2$ and $r_2 \times n_3 \times 1$. For convenience, we can treat these as a matrix, a three-way tensor and a matrix by ignoring dimensions of size one. This leads to the following expression of equation 3.3.2 where \mathcal{W} is the three way tensor in its decomposed form:

$$W_{ijk} = \mathbf{A}_i \cdot \mathbf{B}_{\cdot j} \cdot \mathbf{C}_{\cdot k}.$$

In a manner consistent with the CP-decomposition we denote such a decomposed tensor $\mathcal{W} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{TT}$. It is important to note that the shapes are less consistent: \mathbf{A} is an $n_1 \times r_1$ matrix, \mathbf{B} a $r_1 \times n_2 \times r_2$ tensor and \mathbf{C} an $r_2 \times n_3$ matrix.

As this decomposition only contains contractive products, it is very well expressed as a Tensor Network diagram. Figure 3.4 shows both the general case and the three-way case – the general case shows clearly how we can build a tensor with a large number of dimensions purely out of relatively small three-way tensors.

Bilinear Product

Computing a bilinear product between two vectors and a tensor in the TT-decomposition does not have quite such an efficient form as the CP-decomposition. Primarily this is due to the presence of the three-way tensor \mathbf{B} , which means we are still eventually computing a full tensor product, simply with a smaller tensor. We denote the product as

$$\mathbf{z} = \mathbf{x}^T \mathbf{A} \mathbf{B} \mathbf{C} \mathbf{y}. \quad (3.7)$$

For a specific index this has the form

$$\begin{aligned} z_j &= \mathbf{x}^T \mathbf{A} \mathbf{B}_{\cdot j} \mathbf{C} \mathbf{y} \\ &= \sum_i^{n_1} \sum_k^{n_3} \sum_{\alpha_1}^{r_1} \sum_{\alpha_2}^{r_2} A_{i\alpha_1} B_{\alpha_1 j \alpha_2} C_{\alpha_2 k} x_i y_k. \end{aligned}$$

Eq. (3.7) provides a clear intuition about what the TT-decomposition is doing in the three-way case. Matrices \mathbf{A} and \mathbf{C} project the inputs into a new (potentially smaller) space where the bilinear product is carried out with \mathbf{B} ensuring the result has been pushed back to the appropriate size.

3.3.3 Comparison

In this section we first prove a result on the conditions required for the decompositions to be equivalent. This is followed by a brief discussion of the theoretical differences and similarities

between the two decompositions and whether we can draw any conclusions about their suitability for the task at hand.

It is also worth noting briefly that the gradients are very nearly equivalent. This implies that attempting to learn the parameters directly using gradient descent should work for both or neither, since it seems reasonable to expect similar dynamics.

Equivalence

One condition for the tensor-train to be equivalent to the CP-decomposition is if the central tensor in the TT-decomposition has diagonal, square slices. Intuitively this reduces the tensor product to a matrix product.

Proposition 3.3.2. *The rank R CP-decomposition of a tensor $\mathcal{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{CP}$ is equivalent to a TT-decomposition $[\mathbf{A}', \mathbf{B}', \mathbf{C}]_{TT}$ with both ranks equal to R and $\mathbf{B}'_{ijk} = 0$ except where $i = k$.*

Proof. Consider the slices of $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ formed by fixing the second index and allowing the first and third to vary. If $\mathcal{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_{CP}$ then these slices can be expressed concisely as

$$\mathbf{X}_{\cdot j \cdot} = \mathbf{A} \text{diag}(\mathbf{B}_{j \cdot}) \mathbf{C}^T$$

where $\text{diag}(\mathbf{v})$ denotes the matrix with vector \mathbf{v} along the leading diagonal and zero elsewhere. In the above the vector is in fact the j -th row of the factor matrix \mathbf{B} . It is clear the diagonal matrix must then have shape $R \times R$ so the result has the expected shape $n_1 \times n_3$. We also verify this gives us the appropriate expression for a single element of the tensor:

$$\begin{aligned} X_{ijk} &= \left(\mathbf{A} \text{diag}(\mathbf{B}_{j \cdot}) \mathbf{C}^T \right)_{ik} \\ &= \sum_{\alpha=1}^R A_{i\alpha} \cdot \left(\sum_{\beta=1}^R \text{diag}(\mathbf{B}_{j \cdot})_{\alpha\beta} C_{k\beta} \right) \\ &= \sum_{\alpha=1}^R A_{i\alpha} B_{j\alpha} C_{k\alpha} \end{aligned}$$

where the last step follows from diagonality.

We now construct a TT-decomposition $\mathcal{X} = [\mathbf{A}', \mathbf{B}', \mathbf{C}]_{TT}$ of the same tensor. The ranks of the decomposition will be $r_1 = r_2 = R$ for R the rank of the corresponding CP-decomposition. Let $\mathbf{A}' = \mathbf{A}$ and $\mathbf{C}' = \mathbf{C}^T$. Construct tensor \mathbf{B}' by

$$B'_{ijk} = \begin{cases} B_{ij} & \text{if } i = k \\ 0 & \text{otherwise.} \end{cases}$$

An expression for the central slices of \mathcal{X} in terms of its TT-decomposition which follows from the element-wise definition is

$$\begin{aligned} \mathbf{X}_{\cdot j \cdot} &= \mathbf{A}' \mathbf{B}'_{\cdot j \cdot} \mathbf{C}' \\ &= \mathbf{A} \text{diag}(\mathbf{B}_{j \cdot}) \mathbf{C}^T \end{aligned}$$

by construction of \mathbf{B}' . □

Space Requirements

The number of stored coefficients for a rank R CP-decomposed tensor of size $I \times J \times K$ will be $RI + RJ + RK$. Explicitly storing the tensor would require IJK numbers to be stored. To illustrate the significance of this, consider the case when the tensor being decomposed has all dimensions and rank of equal size. Denoting the size as I , then the explicit tensor will have an I^3 storage requirement while the decomposed tensor needs only $3RI$, a remarkable reduction.

For a TT decomposition with ranks (R_1, R_2) and dimensions $I \times J \times K$ we need $R_1I + R_1JR_2 + R_2K$ storage. If both ranks are equal, this becomes $RI + R^2J + RK$, which is quadratic in the rank. As the ranks have to be integers, this gives us less ability to fine-tune the number of parameters in the decomposition as changing the ranks by small amounts may have large implications on the size of the decomposition.

The CP decomposition is clearly the most appealing in this regard. Not only does it only require a single hyper-parameter to control the size, the fact that it grows linearly in that hyper-parameter makes it much more amenable to adjustment.

3.4 Learning decompositions by gradient descent

We now present some experimental results, with two aims in mind. Firstly, it is instructive to compare the performance of the two tensor decompositions when learnt by gradient descent on some straightforward tasks. Secondly we wish to verify that a tensor decomposition is capable of learning complicated structure by applying it to a classic benchmark and a real world example.

3.4.1 Gradients

It makes sense to check that the gradients of the parameters are well-behaved with respect to the output of the bilinear product. We break the gradient into three parts, one for each of the parameter matrices. As the end goal is the gradient of a vector with respect to a matrix it should naturally be represented by a three-way tensor. Many authors avoid this by vectorising the matrix [60], but for the purposes below it is sufficient to note that the gradients are highly structured, and we can write a simple form for each element one at a time.

CP-decomposition

Let \mathbf{z} be defined as above. The gradient of \mathbf{z} with respect to \mathbf{A} has entries of the form:

$$\begin{aligned} \frac{\partial z_j}{\partial A_{lm}} &= \sum_k^{n_3} B_{lj} C_{lk} x_m y_k \\ &= B_{lj} x_m \cdot \langle \mathbf{C}_{l\cdot}, \mathbf{y} \rangle. \end{aligned}$$

The gradients with respect to \mathbf{C} have the same form:

$$\begin{aligned} \frac{\partial z_j}{\partial C_{lm}} &= y_m \cdot \sum_i^{n_1} A_{li} x_i \\ &= B_{lj} y_m \cdot \langle \mathbf{A}_{l\cdot}, \mathbf{x} \rangle. \end{aligned}$$

The gradient of the final parameter matrix \mathbf{B} has entries

$$\frac{\partial z_j}{\partial B_{lm}} = \sum_i^{n_1} \sum_k^{n_3} A_{li} C_{lk} x_i y_k. \quad (3.8)$$

Curiously, the j and m indices do not appear in the right hand side of eq. (3.8). This appears to suggest that \mathbf{B} may learn a significantly redundant structure. However, during gradient descent this will be multiplied by the back-propagated loss, will alter the update.

A final point to make about these gradients is the degree of multiplicative interactions present. This has been noted to occasionally cause slightly problematic dynamics during gradient descent. The simplest example is to consider gradient descent on the product of two variables: $c = ab$. If a is very small and b is very large then the update applied to b will be very small and the update applied to a will be proportionally very large. Sutskever [90] uses this as motivation for using second-order information during optimisation of RNNs containing similar structures. We prefer to counsel patience – after the above step, the gradients will become perfectly aligned to the scale of the parameters and subsequent updates should be better aligned.

Tensor Train

The gradient of \mathbf{z} with respect to \mathbf{A} has entries of the form:

$$\begin{aligned}\frac{\partial z_j}{\partial A_{lm}} &= \sum_k^{n_3} \sum_{\alpha_2}^{r_2} B_{mj\alpha_2} C_{\alpha_2 k} x_l y_k \\ &= x_l \cdot \left(\mathbf{B}_{mj}^\top \cdot \mathbf{C} \mathbf{y} \right).\end{aligned}$$

The components of the gradient with respect to \mathbf{C} have a similar form:

$$\begin{aligned}\frac{\partial z_j}{\partial C_{lm}} &= \sum_i^{n_1} \sum_{\alpha_1}^{r_1} A_{i\alpha_1} B_{\alpha_1 jl} x_i y_m \\ &= y_m \cdot \left(\mathbf{B}_{jl}^\top \mathbf{A} \mathbf{x} \right).\end{aligned}$$

While the gradient of the elements of \mathbf{B} behave similarly to the CP-decomposition:

$$\frac{\partial z_j}{\partial B_{lmn}} = \sum_i^{n_1} \sum_k^{n_3} A_{il} C_{nk} x_i y_k.$$

This is very similar to the CP-decomposition, very little further can be added. Again the central object – which is in this case a three-way tensor – has highly redundant gradients, although it is not clear what effect this may have on learning. Finally, there is slightly more summation in the tensor-train, but it does not appear to be a dramatic enough recasting to avoid any potential instabilities due to the proliferation of multiplicative dynamics in the gradients.

3.4.2 Random Bilinear Products

Goal

The aim for this experiment is to compare the two decompositions. There are two scenarios to investigate: learning a decomposed tensor when the tensor has the same structure as the data and learning a decomposed tensor when it is incapable of representing exactly the underlying structure.

Experiment Details

The first test uses the following procedure: generate a fixed random tensor \mathcal{T} in the chosen decomposition with rank r_T . Then generate a second tensor \mathcal{W} with rank r_W . At each stage t generate a pair of random input vectors \mathbf{x}_t and \mathbf{y}_t and compute the bilinear products $\mathbf{z}_t = \mathbf{x}_t^\top \mathcal{T} \mathbf{y}_t$ and $\hat{\mathbf{z}}_t = \mathbf{x}_t^\top \mathcal{W} \mathbf{y}_t$. Finally update the parameters of \mathcal{W} to minimise the mean squared error between \mathbf{z}_t and $\hat{\mathbf{z}}_t$ using stochastic gradient descent. In this way we hope to learn an approximation to \mathcal{T} in a similar setting to what might be found inside a neural network.

In practice it is computationally expedient to deal with more than one vector at a time, in all of the results for this section we use a batch size of 32 and a fairly aggressive learning rate of 0.1. Training continues for 250,000 parameter updates, although most of the tensors have stopped making significant improvements by that time. In all tests the input vectors had elements drawn from a uniform distribution over the range $[-1, 1]$ while unless otherwise noted the elements of the decomposition are drawn from a normal distribution with mean 0 and standard deviation 0.1.

Results

The results of attempting to approximate a rank 100 tensor with dimensions $100 \times 100 \times 100$ are shown in figure 3.5 and are more or less as anticipated when attempting to approximate a CP-decomposition with another CP-decomposition. When the approximator is a TT-decomposition, it was found to be very difficult to achieve a good approximation – when the learning rate was dropped to account for the apparent high sensitivity of the loss the still failed to approach the error rates of the CP-decomposition. While we might expect the CP-decomposition to better represent another CP-decomposition due to the shared structure, it is nonetheless remarkable how difficult the TT-decomposition was to train.

In general we see that the quality of the approximation drops consistently with the number of parameters in the approximator, which is to be expected. More unexpected is the variance apparent in the training curves of the CP-decomposition – the very low rank approximations seem more sensitive to small changes in their coefficients. This suggests that some care may need to be taken with the learning rates and the initialisation especially in deeper models as when such difficulties will be compounded. It is also useful to note that over-specifying the rank of the decomposition appears to still aid learning.

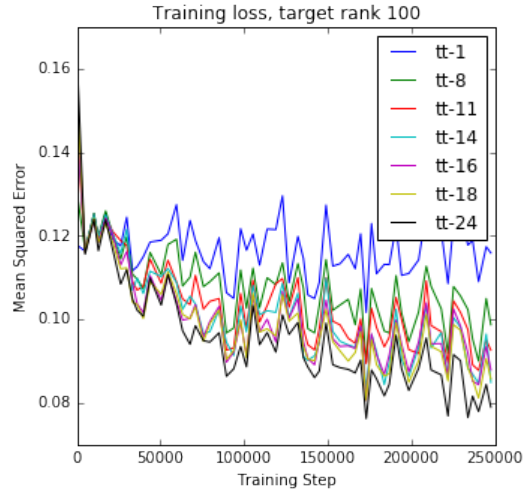
The experiment was repeated with the target tensor expressed as a TT-decomposition with ranks $r_1 = r_2 = 16$ (28,800 coefficients). Results can be seen in figure 3.6. As should be expected, the TT-decomposition performed better. In particular a TT-decomposition matching the rank of the target tensor reached a very low error extraordinarily quickly, performing much better than the CP-decomposition under the same conditions.

These results imply both decompositions are potentially useful. The CP-decomposition learns reasonable approximations even when it does not necessarily reflect the structure of the underlying problem. The TT-decomposition struggles in this situation, but when the problem is structured more favourably it takes greater advantage.

3.4.3 Learning to multiply

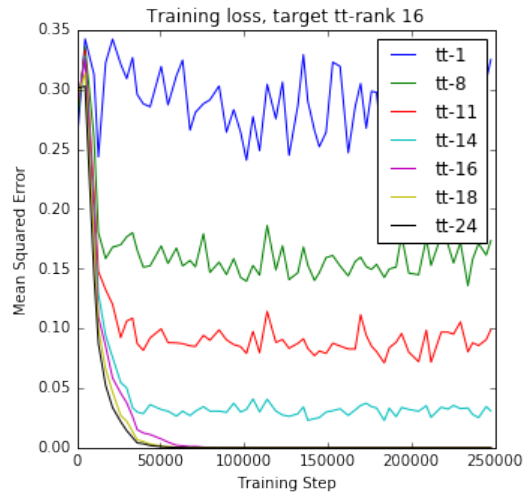
Goal

Learning bilinear functions provides a way to linearly approximate functions of two variables. A motivating example of such a function is element-wise multiplication, also known as the Hadamard product. In this section we investigate briefly the ability of the two decompositions



(a) Attempting to approximate by directly learning factor matrices representing a CP-decomposition of various ranks (b) Attempting to approximate by directly learning a TT-decomposition of various ranks

Figure 3.5: Results for learning direct approximations of a CP-rank 100 tensor.



(a) Attempting to approximate by directly learning a CP-decomposition of various ranks (b) Attempting to approximate by directly learning a TT-decomposition of various ranks

Figure 3.6: Results for learning direct approximations of a tt-rank 16 tensor.

to model such a product. We are particularly concerned with this as generalising this operation provides significant motivation for the current investigation.

Experiment Details

Learning exact element-wise products is a special case and can clearly be done with the CP-decomposition without sacrificing parameter reduction. The question that remains is how closely a given decomposition can approximate an element-wise product, especially if there is some kind of helpful latent structure in the inputs to exploit. We test this by inducing a random structure on the inputs.

This is performed identically to the experiments in section 3.4.2 except that rather than generating a random tensor to provide the targets for training, we simply multiply element-wise the input vectors. We also simplify slightly by drawing the input vectors from $\{0, 1\}$ with equal probability (element-wise multiply in this sense corresponds to a bitwise AND). Again the squared error loss is used, although a momentum term (with coefficient 0.9) was found to greatly assist learning. For a further test, the target output has a fixed permutation applied which removes the diagonal structure from the target but should still be a straightforward task. For this task we would expect that the CP-decomposition drastically outperforms the TT-decomposition.

For a more realistic test we introduce a random structure to the input vectors. This is done by creating smaller vectors and expanding them to the appropriate size by copying the elements to random positions. The mapping of positions is fixed for the duration of a training run. We then experiment with varying the size of these smaller vectors and hence varying the degree of correlation while keeping the rank of the decomposition fixed. Here we have fewer preconceptions regarding the relative performance of the decompositions.

Results

The results, figures 3.7 and 3.8, affirm our expectations. The CP-decomposition is able to consistently reduce error as rank increases while the TT-decomposition appears to struggle.

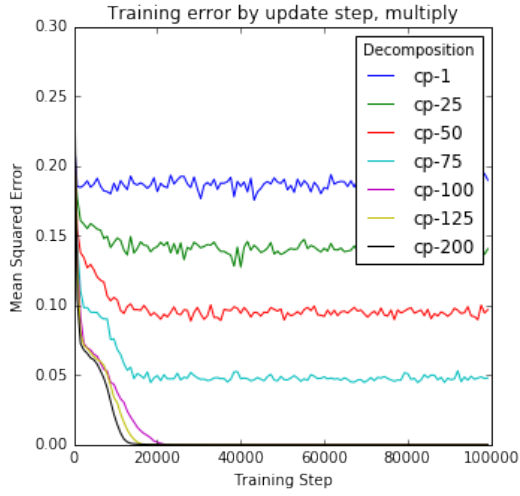
What is more surprising is that even when there is significant structure to the inputs, the TT decomposition is notably worse. One hypothesis is that the CP decomposition simply represents higher rank tensors with the same number of parameters and that this is what allows to better represent more complex structures. This is shown in the other results – the rank 1 decompositions achieve similar performance regardless of the decomposition.

The results again show the CP-decomposition outperforming the TT-decomposition dramatically. In this case this is slightly more surprising as there is more likely to be solutions obtainable in the TT-decomposition.

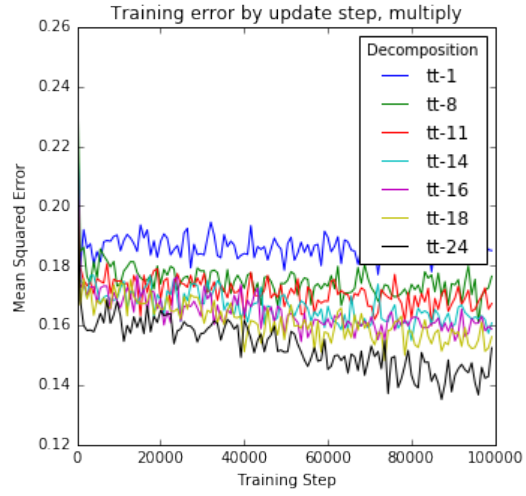
3.4.4 XOR

Goal

Theoretical results suggest that a single tensor layer should be able to solve the exclusive-or problem. This is known to be impossible with a single normal feed-forward layer, so if the tensor layer is successful in this task it would provide a useful demonstration of the difference in expressive power.

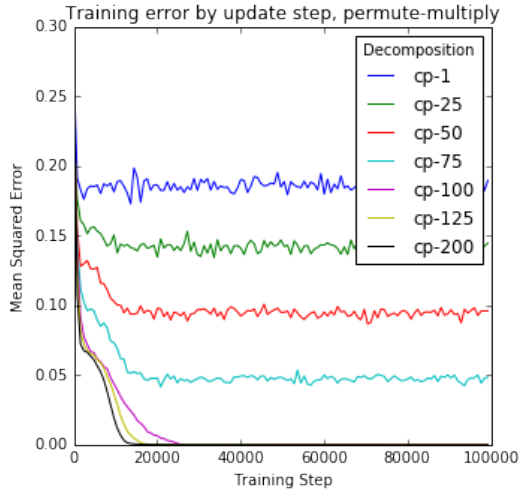


(a) Learning curves for CP-decompositions learning element-wise multiplication

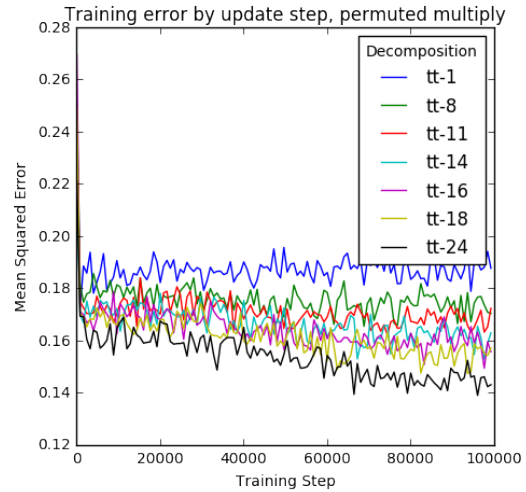


(b) Learning curves for TT-decompositions learning element-wise multiplication

Figure 3.7: Training error for various rank decompositions on the element-wise multiplication task.

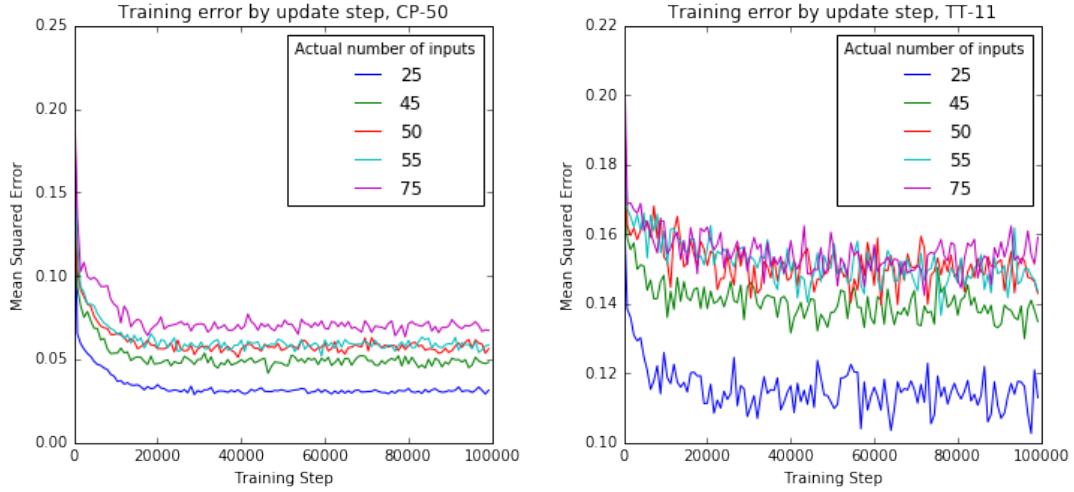


(a) Learning curves for CP-decompositions learning permuted element-wise multiplication



(b) Learning curves for TT-decompositions learning permuted element-wise multiplication

Figure 3.8: Training error for various rank decompositions on the permuted element-wise multiplication task.



(a) Learning curves for a rank 50 CP-decomposition. (b) Learning curves for a rank 11 TT-decomposition.

Figure 3.9: Training error for two decompositions learning element-wise multiplication with varying amounts of structure.

Experiment Details

The truth table for exclusive-or operation can be found in table 3.2. True values were represented as 1 and false as 0. As there are only four data points, all were evaluated for each parameter update. All models were trained using standard gradient descent with a learning rate of 0.5.

Three models were evaluated, a perceptron, a multi-layer perceptron (MLP) and a single layer tensor. All models used a sigmoid nonlinearity on the output, the MLP had 4 hidden units also with sigmoids and the tensor layer used a rank 1 CP decomposition. Models were trained to minimise the cross-entropy, averaged across the data points with training continuing for 1000 epochs. Although not all MLPs had converged by this time, it was sufficient to see the difference between the architectures.

Results

The baseline cross entropy is 0.6931, corresponding to outputting 0.5 for all inputs (or guessing uniformly at random). As expected, the perceptron converges rapidly to this point. Both the MLP and the tensor are able to solve the task, although it is interesting to note that the tensor does so far more rapidly and reliably than the MLP. This is not surprising in light of theorem 1 as this is precisely the kind of relationship we would expect the tensor product to represent naturally.

3.4.5 Separating Style and Content

Goal

In this section we apply the tensor decompositions to a less trivial task. In particular, we address the *extrapolation* problem presented by Tenenbaum and Freeman. [97] The focus in this task is on data that is naturally presented as a function of two factors, termed *style* and *content*. Tenenbaum and Freeman propose the use of bilinear models for such data as

Input 1	Input 2	Out- put
T	T	F
T	F	T
F	T	T
F	F	F

Table 3.2: Exclusive-or truth table

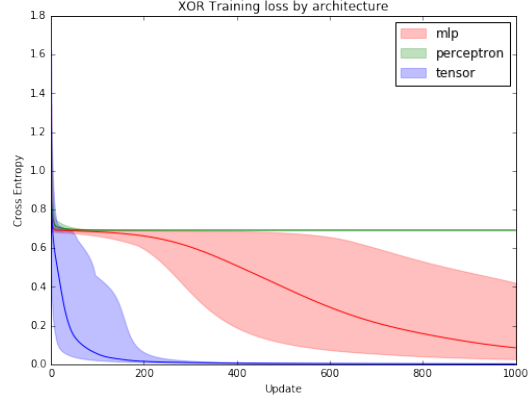


Figure 3.10: Exclusive-or results. Bold line is the mean of 50 runs, shading indicates the minimum and maximum across all runs.

follows. The aim is to verify that we can learn a tensor decomposition of an unknown latent tensor with unknown rank.

Experiment Details

Problem Formulation Let \mathbf{y}_{sc} be an observation vector with style s and content c . It is modelled in bilinear form:

$$\mathbf{y}_{sc} = \mathbf{a}_s^T \mathcal{W} \mathbf{b}_c$$

where \mathbf{a}_s and \mathbf{b}_c are style and content parameter vectors respectively and \mathcal{W} represents a set of basis functions independent of the parameters. A number of tasks involving learning such a model were proposed, but we focus on extrapolation: learning both the parameter vectors and the basis tensor at once in such a way that the model successfully generalises to style and content pairs not seen during training.

We model the style and content parameter vectors as fixed length, dense vectors with a unique vector for each style label and each content label. The parameters in these vectors can be learnt by gradient descent at jointly with the tensor \mathcal{W} . This is equivalent to representing them in a one-of-N encoding (a vector with one value per possible label, all zero apart from the entry corresponding to the label in question) and then multiplying it by a matrix. We refer to such a matrix as an *embedding matrix*. It is clear that if we were to attempt to learn the model without these embedding matrices and simply representing \mathbf{a} and \mathbf{b} with one-of-N encodings that the model would fail completely to generalise. A bilinear product with two such vectors would amount to selecting a single fibre from the tensor, so the tensor could only ever hope to directly store the training data. Introducing a decomposition changes this by forcing the elements of the tensor to depend on one another.

Data One dataset explored in [97] is typographical – this provides a natural source of data defined by two factors: font and character. We refer to the character as the content and the font as the style. We collected a small dataset of five fonts and their italic/oblique variants for a total of ten styles. Variation between styles comes from stroke width, slant and the presence of absence of serifs. From each font the uppercase and lowercase letters of the English alphabet were extracted, totalling 52 different content labels. Data was saved as 32 pixel by 32 pixel greyscale images. To represent images as vectors for the purposes of the above model they are flattened in scanline order: left to right and top to bottom.

Training Details This is a very difficult task as we are expecting the model to produce pictures similar to ones that have never been seen in training. Further, we are trying to train a model with potentially thousands of free parameters based only on a few hundred examples. As this suggests, stopping the model from overfitting was the major challenge. The aim of the experiments was not necessarily to achieve state-of-the-art performance, but to quickly ensure that using a tensor decomposition to model non-trivial interactions was a feasible goal.

To begin we verify that a fairly small model is capable of representing the data. For this we use embeddings of size 25 and a rank 25 CP-decomposition. This model has 78,350 parameters, while a model that contained the explicit $25 \times 1024 \times 25$ tensor would have 641,550. We train the model using ADAM [51], a variant of stochastic gradient descent to minimise the squared error:

$$E = \sum_i^B ||\hat{\mathbf{y}}_{sc} - \mathbf{y}_{sc}||_2^2$$

where B is the number of elements in a batch (26 was used in the following) and $\hat{\mathbf{y}}_{sc} = \mathbf{a}_s^T \mathbf{W} \mathbf{b}_c$ is the predicted image. Both the central tensor and the embedding vectors are updated at every update step. A small amount of l_2 regularisation was found to help prevent overfitting, this involves adding a penalty term to the loss of the form $\lambda ||\mathbf{X}||_F^2 = \lambda \sum_i^m \sum_j^m X_{ij}^2$ for each matrix in the model.

Results

Figure 3.11 provides a visual inspection of what the model was able to learn. Images were generated by finding the style and content labels for a two examples and linearly interpolating first the content vector and then the style vector, generating an image with the intermediate vectors at each step.

The start and end images are not perfect; the model was small and unable to capture the training data exactly. During experimentation larger models were found to fit the data very well, but overfitted rapidly, which is why a smaller model was used for these results.

Intermediate stages of the interpolation have no real meaning but they indicate that the model has not simply learned to recall individual pairs of labels. We see that as either the content or the style shifts from one to another, the elements of the source image which are not shared with the target fade out and are replaced.

Results on unseen data are shown in figure 3.12. These show that the model does appear to capture some of the salient information for separating the style and content. In particular, the general shapes of the letters are roughly appropriate and it has begun to capture some information about the presence or absence of serifs and general slant. Figure 3.13 shows a harder example – the lowercase ‘a’ has considerable variation among various fonts and the model is unable to guess what would be appropriate for an unseen example.

Overall, these results show the tensor was capable of learning a model for the data. This suggests that the aim of inserting such tensor products into more complex neural networks is feasible, reinforcing the promising theoretical results.



(a) Linearly interpolating between two content vectors ('T' to 'm') with style vector fixed.



(b) Linearly interpolating between two style vectors with fixed content.

Figure 3.11: Learned style and content representations.



(a) Actual images from the dataset.



(b) Model output, having never seen these specific pairs during training.

Figure 3.12: Example of images generated from unseen style and content pairs, three letters from three different fonts.



(a) From the data



(b) Generated

Figure 3.13: A difficult example.

Chapter 4

Proposed Architectures

In this chapter we use the intuitions gathered from related works to propose novel classes of architectures employing tensor decompositions to implement multiplicative connections. There are two key principles: simplicity and modularity. We want to design networks that have no extraneous components, such that every element of the network has a clearly defined role.

4.1 Incorporating tensors for expressivity

We have established that incorporating tensors into networks is a promising line of inquiry, all that remains is how precisely to formulate it.

4.1.1 Biases

Neural networks typically require biases to be able to express a full range of transformations. We would expect a tensor layer to be no different. A common way to conceptualise the addition of a bias for a given layer is to incorporate it into the weight matrix by adding an additional row and inserting a corresponding input which has its value fixed to one. With a bilinear product gives more than just the addition of a bias vector. We use $\tilde{\cdot}$ to represent altering the inputs and weights in this way: $\tilde{\mathbf{x}}$ is an input vector with an additional 1 appended and $\tilde{\mathbf{U}}$ is a matrix with a corresponding additional row.

Generalising this construction to a three-way tensor we have to take an $n_1 \times n_2 \times n_3$ tensor to a $n_1 + 1 \times n_2 \times n_3 + 1$. This leads to the addition of a vector as well as two more matrix products. For a three-way tensor \mathcal{W} ,

$$\tilde{\mathbf{x}}^T \tilde{\mathcal{W}} \tilde{\mathbf{y}} = \mathbf{x}^T \mathcal{W} \mathbf{y} + \mathbf{U} \mathbf{y} + \mathbf{V} \mathbf{x} + \mathbf{b}. \quad (4.1)$$

This is very similar to the manner in which the states are computed in the vanilla RNN, just with the addition of the multiplicative tensor interactions.

Applying the same process when the tensor is represented in the CP decomposition does not have the same result. Instead it simply results in adding bias vectors to two of the internal matrix products. Let $\mathcal{W} = [A, B, C]_{CP}$, then simply appending ones to the inputs results in

$$\tilde{\mathbf{x}}^T \tilde{\mathcal{W}} \tilde{\mathbf{y}} = \mathbf{B}^T ((\mathbf{A} \mathbf{x} + \mathbf{b}_x) \odot (\mathbf{C} \mathbf{y} + \mathbf{b}_y)).$$

This is quite different to equation (4.1) meaning that if we wish to insert such a tensor product into a neural network we must choose whether to leave the biases in the decomposition or treat the bias matrices separately. The latter provides less control over the parameters, as there are now two matrices unaffected by the rank, but it might be helpful in that we could use a very low rank decomposition and still maintain a baseline behaviour.

4.1.2 Generalised Multiplicative RNN

The simplest architecture way to incorporate a decomposed tensor is to use it to generalise the Multiplicative RNN and Multiplicative Integration RNN [91, 104]. To do this, we simply replace the various linear operations with a bilinear form with appropriate biases:

$$\mathbf{h}_t = \tau \left(\tilde{\mathbf{x}}_t^\top \tilde{\mathbf{W}} \tilde{\mathbf{h}}_{t-1} \right)$$

Choosing to keep the biases elements separate from the decomposition gives a form which captures vanilla RNNs as well as several types of multiplicative RNNs:

$$\mathbf{h}_t = \tau \left(\mathbf{x}_t^\top \mathbf{W} \mathbf{h}_{t-1} + \mathbf{U} \mathbf{h}_{t-1} + \mathbf{V} \mathbf{x}_t + \mathbf{b} \right).$$

We term this the Generalised Multiplicative RNN (GMR). Unfortunately this network will still exhibit the same vanishing gradients as the vanilla RNN as the state updates still pass through a matrix (albeit one modulated by the new input).

4.2 Gates and Long Time Dependencies

To address these vanishing gradients, we turn to the same solution as the LSTM and GRU: an additive state update. A naive purely additive state update would solve the issues with the gradient vanishing, but causes a number of other issues. A better method is to use a *gated* recurrence. The LSTM and the GRU both use such a method to produce new states, but they differ slightly. We investigate the implications of both schemes and decide that the scheme employed in the GRU has appealing properties which allow us to design simple yet flexible architectures with well defined components.

4.2.1 Naive Addition

The simplest possible approach is to have the network compute a candidate vector of hidden states \mathbf{z}_t and compute new hidden states as

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \mathbf{z}_t. \tag{4.2}$$

Recalling equation (2.3), the problematic term was $\nabla_{\mathbf{h}_{k-1}} \mathbf{h}_k$, the gradient of a hidden state with respect to its immediate precursor. If the state is computed by equation (4.2):

$$\nabla_{\mathbf{h}_{k-1}} \mathbf{h}_k = \mathbf{I} + \nabla_{\mathbf{h}_{k-1}} \mathbf{z}_k.$$

Adding the identity matrix ensures the eigenvalues of the gradient are at least one, which negates the potential to vanish.

While the gradients will not vanish, they may still explode. Pascanu et al. show that a necessary condition for exploding gradients, when using the hyperbolic tangent activation, is that the largest eigenvalue of the recurrent weight matrix is greater than one [77]. This is shown by observing that the l_2 norm of the gradient¹ is upper bounded by the product of the norms of the recurrent weight matrices and the gradient of the nonlinearity. For our purposes it suffices to observe the eigenvalues of the gradient with this additive scheme are always at least one, so unless the other component is zero, the gradients must necessarily explode.

¹ The l_2 norm of a matrix is the norm induced by the l_2 norm on vectors:
 $\|\mathbf{A}\|_2 = \sup \{ \|\mathbf{Ax}\|_2 : \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2 = 1 \}$ which turns out to be the square root of the largest eigenvalue.

With care and gradient clipping [77] exploding gradients can be managed, but there is a more intuitive reason why this network is naive: a purely additive network will struggle to forget. At each step, it will always add to the current state. Therefore, the candidate state must contain negative values to ever remove any information from the state. This suggests two things: the candidate state must depend on the previous state (so that it knows what to remove) and that the network will need to be very precise with its scale so that it can actually remove information that is no longer necessary. This last point suggests that the network will have trouble if the non-linearity applied to the candidate states is bounded – if it erroneously adds during training it will struggle to learn to subtract. An unbounded non-linearity would likely exacerbate the gradient issues, which makes this sort of a network challenging to design.

4.2.2 Gates

LSTMs and GRUs manage their additive connections by applying multiplicative gates. We will therefore investigate the solutions they present, with a view to choosing the best. We refer to the LSTM’s gate as a *forget* gate, which is consistent with the literature. This type performs best in the presence of other gates [33, 48], which is explained by our theoretical analysis. The gate on the GRU is slightly different and has a number of advantageous properties. We term this the *convex* gate and find it is sufficiently flexible to allow us to remove the dependency of the candidate state on the hidden state, which is desirable in terms of simplifying the architecture as well as controlling the gradients.

Before delving into their specific properties, we first describe the gating mechanisms. Denote by p_t the gate signal at time t (which is a function of both the input and the current state and bounded in $[0, 1]$). We also use h_t to refer to the hidden state at time t , the object of interest, and z_t will refer to the candidate update computed by the rest of the network. In general these elements will be vectors, but in the following we assume scalars with no loss of generality as the vector forms contain only element-wise operations.

We can define the forget gate’s recurrence as

$$h_t = p_t h_{t-1} + z_t$$

for z_t the new candidate state which is typically a function of both h_{t-1} and the input at time t . This allows the network to wholly replace the state with a new value, but as will be shown in detail it has no control over the incoming flow of information. For exact equivalence with the LSTM, it suffices to note that the z_t could be computed in any way, including with its own gating mechanisms.

The convex gate makes the following modification:

$$h_t = p_t h_{t-1} + (1 - p_t) z_t.$$

This allows the one gate signal to control the acceptance or rejection of new information.

In order to analyse the possible behaviours of these gates, we observe that each state h_t can be expanded as a weighted sum over all z_i for $i < t$. In this sense, we can think of the states as being a sliding-window sum over candidates where the shape of the window is defined recursively by the p_i values.² With this view in mind, we note that we can think of the gate signals p_i as providing an *attention* mechanism which controls how the network attends to the information from past time steps.

²There is in fact no sliding window – the sum is always over all values. However given the range of window shapes the gating mechanism is capable of producing it is reasonable to think about it as something close to a one-dimensional convolution.

Forget Gate

The forget gate is designed to allow the network to completely replace a state, hence its name. Clearly for this to happen the gate signal simply needs to be zero. Further, the default behaviour is to pay more attention to recently acquired information. If we define the initial state $h_0 = 0$, we can rewrite the first few states as

$$\begin{aligned}h_1 &= z_1 \\h_2 &= p_2 z_1 + z_2 \\h_3 &= p_3 h_2 + z_3 \\&= p_3 (p_2 z_1 + z_2) + z_3 \\&= p_3 p_2 z_1 + p_3 z_2 + z_3\end{aligned}$$

In general,

$$h_t = \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t p_j \right) \cdot z_i + z_t.$$

Each state is a weighted sum of the past states, with strictly non-increasing weights as we go back in time. Note that the coefficient of each past candidate is the product over all gate signals from the current time back to the candidate in question. As the gate signals are in $[0, 1]$, this is likely to correspond to a very steep decay.

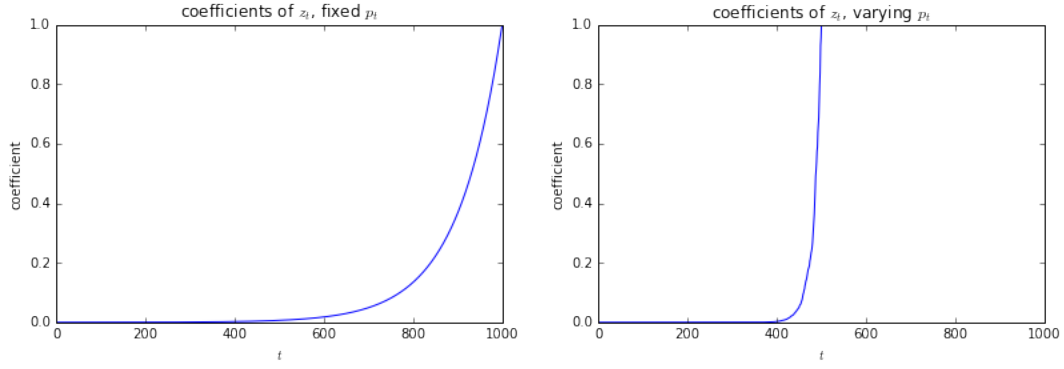
Figure 4.1 shows some consequences of this scheme. It is clear that it offers a blunt choice – exponential decay over time or accumulating all states from a certain point nearly equally. This means that the mechanism for producing candidate states is very important.

Consider the case where the network needs to remember an event from one time step some time in the past and ignore all further input. Achieving this with the forget gate can be done in two ways. Firstly, the gate values can be clamped to one and the candidates to zero. Clamping the gates to one produces a rectangular window; the state will be an evenly weighted sum of all candidates from the point the gate switches on. If the goal is to pick out a single candidate and remember it for a long time, all subsequent candidates will have to be zero. The second method of remembering a single candidate state is by keeping the exponential window and simply producing the same candidate at every future time step (with appropriate scaling).

Both of these require a significant degree of co-operation between the gate and the mechanism that produces candidate activations. Further, the production of candidate activations has to be quite flexible and most likely depend on the previous state. The LSTM achieves this with further gating of inputs and outputs, enabling the candidate to depend heavily on previous states. We hypothesise that having to learn this degree of co-operation hinders learning by gradient descent as it requires a carefully tuned balance to be maintained through the training process or gradients will vanish or explode.

The shape of this window affects access to past information which can be an issue during training. Intuitively it would be good for the information stored in the states to persist so that its presence will affect the gradients being back-propagated which will affirm or suppress its presence. This is going to be an issue when the newly accumulated information decays away exponentially in time, as will be the case early on when the gate values will be essentially random. A solution to this issue commonly used in practice is to ensure the bias of the unit computing the gate value is initialised to a high positive value to force it to stay early in training [48].

The converse is also an issue. If solving the task at hand requires storing information for a long time, then at some point during training the forget gate will end up in a position where



(a) Coefficients of z_t resulting from a fixed, (b) Coefficients of z_t random p_t up to $t = 500$ high value of all p_t under the forget gate and all p_t after 500 set to 1. scheme.

Figure 4.1: Different window shapes produced by the forget gate. As each coefficient at time t is the product of all gate values from t down to 1, the range of possible shapes is limited.

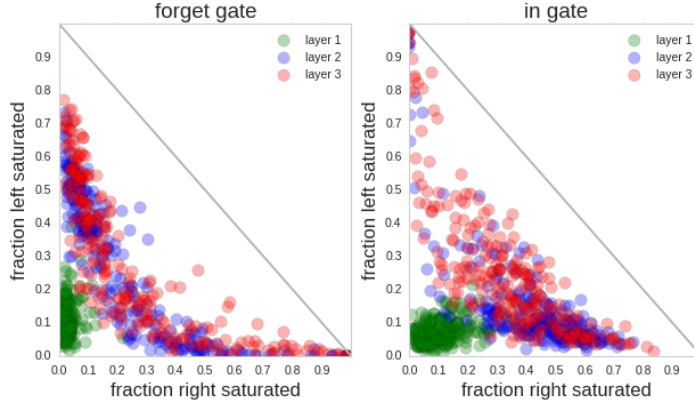


Figure 4.2: Saturation rates of 3-layer LSTM. The input gate modulates the candidate update, while the forget gate is as defined in this chapter. Figure from [50], reproduced with permission.

it is very high for a long time during which the states accumulate in an unbounded fashion. A small change to the state early on this period could trigger exponential growth over time (as the candidate production depends multiplicatively on the previous state).³

Further, we hypothesise that the strong connection required between the candidate production mechanism and the forget gate indicates a strong potential for redundancy. Some experimental work with LSTMs indicates that this may be the case. In [50] the authors trained a large LSTM on a language modelling task and plot the gates by the fraction of time they spend left-saturated (value < 0.1) or right-saturated (value > 0.9). This is reproduced in figure 4.2 and shows that the forget gates in general tend to spend more time in their lower ranges while inputs appear most likely to be in the higher ranges of their activation. In general, it seems to suggest most of the state is replaced at each time step (although looking at the input gate saturations, some cells must almost never change their value). With an LSTM this does not preclude it from storing data for long periods, but it means that in order to do so it must learn an approximate identity mapping, which seems like wasted effort.

It is also worth considering how this gating scheme affects the gradients during back-

³ This is an example of the butterfly effect noted previously.

propagation. In order to update the weights producing the gate signal, we will need the gradient of the hidden state with respect to the gate signal itself at some earlier time $i < t$:

$$\frac{\partial h_t}{\partial p_i} = \left(\prod_{k=i+1}^t \frac{\partial h_k}{\partial h_{k-1}} \right) \frac{\partial h_i}{\partial p_i}$$

by the chain rule. Differentiating equation 4.2.2 appropriately we get

$$\frac{\partial h_t}{\partial p_i} = \left(\prod_{k=i+1}^t p_k \right) h_{i-1}.$$

This makes it clear that the gating does not just happen during the forward pass. The gradients continue to be gated in the same way as the states. This suggests the gate could struggle to learn long time dependencies as information that is erroneously decayed during the forward pass will have a similarly small contribution to the gradients.

Repeating the process with respect to the candidate activations, it is clear that precisely the same argument applies.

$$\begin{aligned} \frac{\partial h_t}{\partial z_i} &= \left(\prod_{k=i+1}^t \frac{\partial h_k}{\partial h_{k-1}} \right) \frac{\partial h_i}{\partial z_i} \\ &= \left(\prod_{k=i+1}^t p_k \right) \cdot 1. \end{aligned}$$

During training this will be multiplied by the back-propagated error. If the gates exhibit a strong decay, then candidate activation will also be forced to update primarily on local information.

Convex Gate

In order to repeat the analysis with the convex gate, we need to more carefully define initial conditions. In particular, let $p_0 = 0$ and correspondingly $h_0 = z_0$ be some initial state.

$$\begin{aligned} h_0 &= z_0 \\ h_1 &= p_1 z_0 + (1 - p_1) z_1 \\ h_2 &= p_2 p_1 z_0 + p_2 (1 - p_1) z_1 + (1 - p_2) z_2 \\ h_3 &= p_3 p_2 p_1 z_0 + p_3 p_2 (1 - p_1) z_1 + p_3 (1 - p_2) z_2 + (1 - p_3) z_3 \end{aligned}$$

giving rise to the general form

$$h_t = \sum_{i=0}^t \left(\prod_{j=i+1}^t p_j \right) (1 - p_i) z_i$$

ensuring we define $\prod_{j=i+1}^t p_j = 1$ if $i = t$. This is very close to what we had above, but the extra $1 - p_i$ in the coefficients has significant impact. Firstly, we can prove the following statement which shows that the sum over states is a convex sum.

Proposition 4.2.1 (Conservation of attention). *At any time $t > 1$,*

$$\sum_{i=0}^t \left(\prod_{j=i+1}^t p_j \right) (1 - p_i) = 1.$$

The coefficients of all previous states sum to one.

Proof. If we expand the brackets, we can produce a telescoping sum. To ensure it telescopes appropriately, we first pull the first and last terms out of the summation.

$$\begin{aligned} & \sum_{i=0}^t \left(\prod_{j=i+1}^t p_j \right) (1 - p_i) \\ &= \prod_{j=1}^t p_j (1 - p_0) + (1 - p_t) + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t p_j \right) (1 - p_i) \\ &= \prod_{j=1}^t p_j (1 - p_0) + (1 - p_t) + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t p_j \right) - \left(\prod_{j=i}^t p_j \right) \\ &= \prod_{j=1}^t p_j (1 - p_0) + (1 - p_t) - \prod_{k=1}^t p_k + p_t \\ &= \prod_{j=1}^t p_j - \prod_{k=1}^t p_k + 1 - p_0 + p_t - p_t \\ &= 1 \end{aligned}$$

□

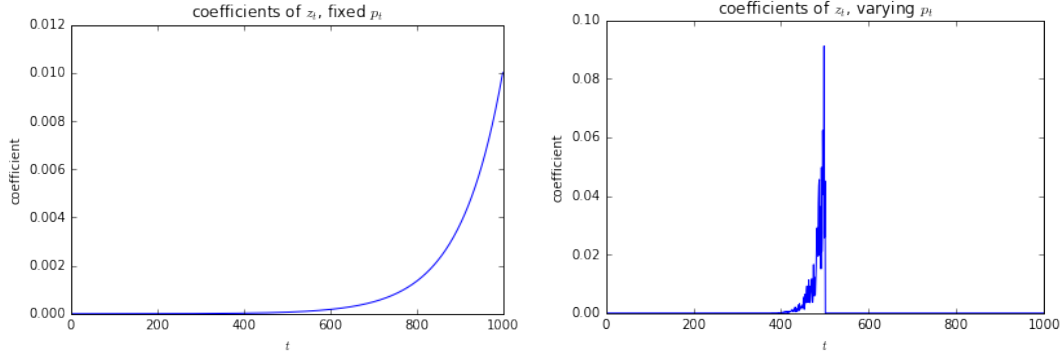
Remarkably, this does not require p_i to be bounded at all. However, in order to compare directly, we consider the common case when the gate signals are between zero and one. The immediate effect of this lemma is that the scale of the states depends only on the scale of the candidate activations. Each state is now a weighted mean of past candidates, rather than a weighted sum.

Figure 4.3 illustrates two of the possible shapes this can take. It is clear that this gating scheme provides a different set of possible shapes to work with. Specifically, if the gate mechanism starts outputting values very close to one, the window applied to past states takes the form of a distinct spike over a small range of activations. This occurs because if p_i is very close to one, $1 - p_i$ is very close to zero, so state h_{i-1} is going to be carried over with only a very minor contribution from z_i . Figure 4.4 further emphasises the difference, making it clear that the convex gate is capable of representing a fundamentally different set of window shapes than the forget gate, despite their similar formulation.

This enhanced range of window shapes makes these gates appealing as they may be able to reduce interdependence between the gate mechanism and the candidate production mechanism. This might allow us to remove the candidate's dependence on the previous state which should help solve the potential instabilities noted above.

The downside of the system is that it may still struggle to learn long time dependencies early. It would be straightforward to initialise the gate so that it has a mean activation of 0.5, but this would correspond to an exponential decay of information. Unlike the forget gate simply increasing the mean activation (for example by initialising the bias to a high positive value) will not correspond to a flatter window.

To finish the analysis we consider again the gradients required for back-propagating error through the gate. The gradient of the hidden state h_t with respect to a previous gate signal p_i



(a) Coefficients of z_t resulting from a fixed, (b) Coefficients of z_t random p_t up to $t = 500$ high value of all p_t under the convex gate and all p_t after 500 set to 1. scheme.

Figure 4.3: Different window shapes produced by the convex gate. A simple adjustment to the formula a much wider range of shapes.

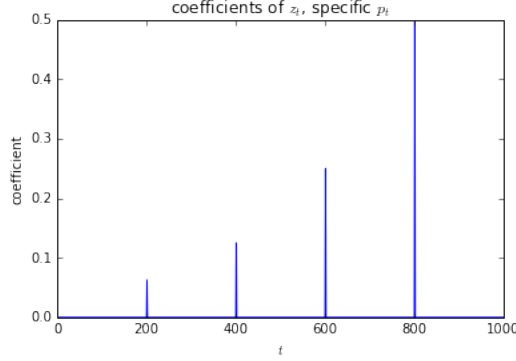


Figure 4.4: Setting specific p_i to 0.5 picks out a set of past candidates, albeit with an exponential decrease in their weighting.

has the form:

$$\begin{aligned} \frac{\partial h_t}{\partial p_i} &= \left(\prod_{k=i+1}^t \frac{\partial h_k}{\partial h_{k-1}} \right) \frac{\partial h_i}{\partial p_i} \\ &= \left(\prod_{k=i+1}^t p_k \right) (h_{i-1} - z_i) \end{aligned} \quad (4.3)$$

While the gradient of the state with respect to a prior candidate is

$$\begin{aligned} \frac{\partial h_t}{\partial z_i} &= \left(\prod_{k=i+1}^t \frac{\partial h_k}{\partial h_{k-1}} \right) \frac{\partial h_i}{\partial z_i} \\ &= \left(\prod_{k=i+1}^t p_k \right) (1 - p_i). \end{aligned} \quad (4.4)$$

The pattern here is fundamentally the same as the forget gate in that the gradient is simply the coefficient assigned to z_i in the weighted sum that produced h_t . Hence all discussion of the differences in the gates' forward behaviour apply equally during the backward pass.

Equation (4.3) is interesting – the gradient of the state with respect to the gate value depends not just on the preceding state but on the difference between the state and the proposed update. This suggests that during training what drives the updates is the possible changes that could be made to the state, rather than just the state values itself. Meanwhile, equation (4.4) shows that the “conservation of attention” in proposition 4.2.1 also applies on the backward pass. This suggests that as long as the scale of the candidate activations are under control there will be no problems with exploding gradients.

Conclusions

The convex gate by itself seems more capable than the forget gate as it is able to represent more useful distributions over past candidate states. However, the combination of forget gate and input gate employed in the LSTM are able to together represent the same behaviours, if not more. Despite this, we find the convex gate more appealing as it provides the opportunity to design clean and modular architectures in which each component has a precise role which does not overlap with others.

4.3 Gated Tensor RNNs

To design a gated RNN, there are three choices to make: how to gate the recurrence, how to compute gate values and how to compute candidate updates. We are now in a position to make these choices on the basis of chapter 3 which discussed how to use tensors to model binary functions and the above section 4.2.2 informs our choice of gate.

4.3.1 Choices

Form of the gate

We now consider how to compute the \mathbf{p}_i values. The role of the gate is to manage the flow of new information into the states. The decisions it has to make are whether to forget what is currently stored and completely overwrite it or reject a candidate update if what is present in the state is somehow more useful. Of the two competing choices, we prefer the convex gate. This is because it is capable of expressing a rich range of distributions over past states in a single operation and without exploding gradients. This gives a recurrence of the form

$$\mathbf{h}_t = \mathbf{p}_t \odot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{p}_t) \odot \mathbf{z}_t. \quad (4.5)$$

Gate calculation

In order to decide whether incoming information should overwrite the current state values, the gate will need to see the current values and the new information. The binary nature of this task suggests the use of a tensor unit as discussed in chapter 3. Secondly if we want the full utility of the convex gate, it seems logical to imbue it with full expressive power – the bilinear tensor is a perfect fit.

It is also necessary to choose a non-linearity. Much of the above analysis assumes the gate signals are in the range $[0, 1]$ although this was not used in the proof of proposition 4.2.1. The typical approach is to ensure this smoothly with a sigmoid. Although sigmoids have undesirable properties when used in feed-forward networks [27] they remain the standard choice when this kind of gating is desired [74, 73] and we find no reason to deviate from this.

The desired gate value is then calculated as:

$$\mathbf{p}_t = \sigma \left(\tilde{\mathbf{x}}_t^\top \tilde{\mathbf{W}} \tilde{\mathbf{h}}_{t-1} \right). \quad (4.6)$$

Following the results in chapter 3 we intend to represent the tensor in a CP decomposition. It remains a choice to decide how to deal with the biases – in the following experiments we tend to include results for both choices for comparison. In particular, keeping the biases implicit potentially provides the ability to regularise the model by restricting the rank, although it needs to be seen empirically whether this affects adversely the representative power of the model.

Using a tensor in this fashion provides an interesting possibility: the ability to write to the memory associatively. One interpretation of the tensor product is that it computes similarity measures between its inputs. Under this interpretation, the gate is constantly computing weighted similarities between the input and the current state. It is possible then for the gate to react if a certain pattern in the input matches a pattern stored in the state and admit or deny a corresponding update. This example shows the power of incorporating multiplicative connections – this mode of behaviour is very difficult to realise when the inputs share only an additive connection.

Candidate update

All that remains is to determine how to derive a candidate state update. In all of the architectures surveyed apart from the Strongly Typed variants [4] it is assumed that this needs to be a binary function of the current input and the previous hidden state. Intuitively this needs to be the case with the LSTM style forget gate, but as we have decided against that it is worth questioning this assumption.

Removing this step makes each role in the network clearly defined. As each input comes in, it is embedded into the state space. The gate decides whether it is important information and the state is updated additively. The candidate production mechanism only needs to behave in as a feature detector, transforming the input into a representation that captures its essential structure in a way amenable to additive composition.

This scheme is appealing due to its simplicity and separation of roles, although two questions remain to be answered: how does this affect the expressive power of the network and how should we compute the necessary representation. The former is likely to be a best answered empirically, although we make a brief note here. One aspect of the GRU and LSTM that is certainly lost is what we term the “carousel” behaviour. A method for these networks to maintain information in their hidden states for long time periods is for the candidate production mechanism to learn an identity mapping from hidden states to hidden states. This behaviour is impossible to realise unless the candidate state depends on the previous hidden state. Note that this is also the only way that the Vanilla RNN can remember information for long time periods. Given the failures of the Vanilla RNN, we hypothesise that removing this mode of behaviour should be no great loss. Further, we suggest it may help optimisation if there is only one clear way in which to store information. This leads us to conclude that we want to remove the hidden state from the candidate calculation.

What remains is to determine how to compute the candidate from the input. We wish to model a unary mapping from the input to the candidate, so a logical way to proceed is with a standard feed-forward layer:

$$\mathbf{z}_t = f(\mathbf{W}_{in} \mathbf{x}_t + \mathbf{b}_{in}). \quad (4.7)$$

We experimented with a number of options for the non-linearity f and found the best performance was consistently achieved with either a linear rectifier $\rho(x) = \max(0, x)$ or no

non-linearity at all. This echoes trends in feed-forward networks away from saturating non-linearities towards unbounded, piecewise linear activation functions [28, 37]. It is interesting that for many tasks a completely linear candidate performed optimally – the only hypothesis is that the sigmoid applied to the gate is enough non-linearity to preserve the representative power of the network as a whole and that a linear layer will have very strong gradients and potentially preferable training dynamics [84].

While a single feed-forward layer was sufficient in our experiments, there is scope to expand this, potentially in a task-dependent manner. For example this could be a deeper feed-forward architecture, a convolutional neural network if the inputs have appropriate structure or even an identity map if the input to the network is already high-level features.

4.3.2 Tensor Gate Unit

Bringing together equations (4.5), (4.6) and (4.7) gives us the proposed Tensor Gate Unit (TGU) which computes its hidden states as follows:

$$\begin{aligned}\mathbf{h}_t &= \mathbf{p}_t \odot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{p}_t) \mathbf{z}_t \\ \mathbf{p}_t &= \sigma \left(\tilde{\mathbf{x}}_t^\top \tilde{\mathcal{W}} \tilde{\mathbf{h}}_{t-1} \right) \\ \mathbf{z}_t &= f(\mathbf{W}_{in} \mathbf{x}_t + \mathbf{b}_{in})\end{aligned}$$

We will represent the tensor with a CP decomposition. We still have to decide whether to incorporate the biases in the decomposition or keep them separate. A straightforward inclination is to keep them separate for the reasons outlined in section 4.1.2, leading to the following gate equation:

$$\begin{aligned}\mathbf{d}_t &= \mathbf{B}^\top (\mathbf{A} \mathbf{x}_t \odot \mathbf{C} \mathbf{h}_{t-1}) \\ \mathbf{p}_t &= \sigma(\mathbf{d}_t + \mathbf{U} \mathbf{h}_{t-1} + \mathbf{V} \mathbf{x}_t + \mathbf{b}).\end{aligned}$$

We also consider the alternative where the bias matrices are folded into the decomposition:

$$\begin{aligned}\mathbf{d}_t &= \mathbf{B}^\top ((\mathbf{A} \mathbf{x}_t + \mathbf{b}_x) \odot (\mathbf{C} \mathbf{h}_{t-1} + \mathbf{b}_h)) \\ \mathbf{p}_t &= \sigma(\mathbf{d}_t)\end{aligned}$$

Chapter 5

Evaluation of architectures

In this chapter we present experimental results comparing the performance of the proposed architectures. something something something.

5.1 Synthetic Tasks

We first evaluate the performance of the architectures on some synthetic tasks. These tasks are designed to be difficult and to exercise various faculties of the model.

5.1.1 Addition

Task

This task is designed to test the networks' ability to store information for long time periods. It first featured in [44], although we use the slightly different formulation found more recently in [56]. The problem is a common benchmark for RNNs and has featured in a number of recent works ([3, 38, 5, 71] for example).

The inputs for this task are sequences of length T . There are two inputs, one sampled from a uniform distribution over the range $[0, 1]$ while the other is zero except for two locations where it is one. The location of the first one is always chosen to be earlier than $T/2$ while the second is after. The goal is to output at the last time step the sum of the two random values that were presented when the second input was one. Pseudocode for an algorithm to generate sequences is presented in the appendix, section B.1. This task becomes harder as the T increases because the length of time the numbers have to be remembered for increases.

Experiment Setup

We test a number of architectures on sequences with T of 250, 500 and 750. As we are concerned purely with the ability of the networks to solve the task, we generate a new batch of sequences at each step. All models were allowed 8 hidden units and trained on batches of 8 data items to minimise the squared error between the output of the network at the final time step and the target value. The TGU was allowed a rank 4 tensor, giving it the smallest number of parameters of the gated architectures tested. Updates were applied using ADAM [51] with a base learning rate chosen from $\{0.1, 0.01, 0.001, 0.0001\}$ according to speed of convergence. For all runs, training continued for 1800 updates. While some authors have found they require significantly more updates for success with certain architectures (up to 10,000,000 in the case of [56]) we found most successful architectures converge within 500 on most sequences.

While previously reported results for this task stop at $T = 750$, the most successful architectures still converge rapidly to a solution at this stage. In order to more fully test their limits we increase T to 1,000, 5,000 and 10,000. The same parameters were used but only the architectures tested that were able to solve the task at $T = 750$ were applied to these much more challenging tasks. In order to reliably find solutions at length 10,000 we had to increase the number of hidden units to 32, the rank of the tensor decomposition was increased proportionally to 16.

The models tested were the LSTM, GRU, IRNN [56] and our TGU. We use a standard LSTM without peephole connections. For the TGU we separate the bias matrices and use a layer of ReLUs for the candidate production. All models use a linear output layer to reduce the hidden states back to a scalar.

Results

The baseline result for this task is a mean squared error of 0.1767 which corresponds to predicting 1 for every input. Many architectures are unable to perform significantly better than this. Figure 5.1 shows that we were unable to find a solution with the IRNN or LSTM within the time allotted, while the GRU and TGU consistently did so very rapidly – in less than 1,000 iterations which several orders of magnitude faster than previously published results [56, 3, 38].

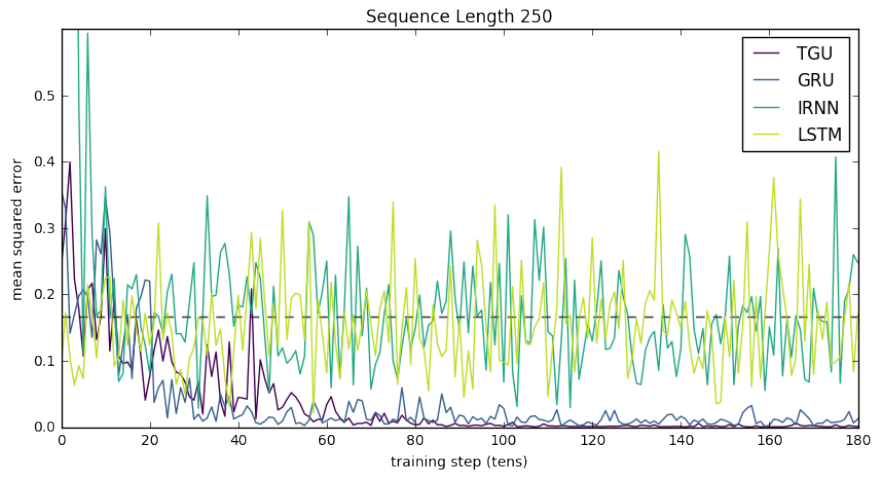
The figures also show that this task can be somewhat frustrating. Due to the online nature, the loss curves are very noisy, so it is almost impossible to determine whether a model is making progress. Further, in accordance with a number of others’ results [56, 3] that many architectures will sit on or around the baseline showing no signs of progress for thousands of training steps before suddenly and rapidly converging to a solution. This behaviour is not particularly evident in the presented results (although the GRU in figure 5.1c is close) as we ceased training fairly early compared to many authors as our successful architectures converged very rapidly.

On sequences up 750 steps, there is little to separate the GRU and the TGU. As they share a gating mechanism, it is probable that this is a key factor in their success. As discussed in chapter 4 this convex gating mechanism is naturally capable of representing exactly the behaviour required to solve this task – remembering precisely a small subset of inputs.

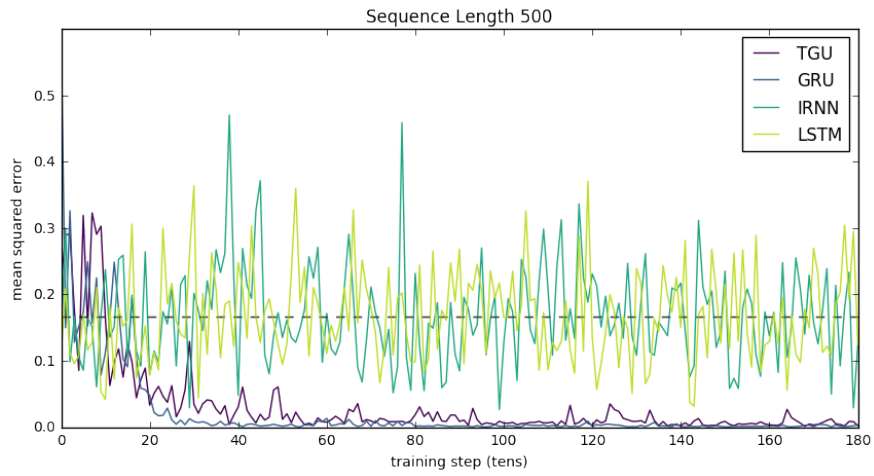
We also found that the TGU performed much better on this task with a ReLU layer producing candidate updates, which is contrary to findings on several later tasks in which linear candidates were preferred. A possible reason for this is that it is very easy for the ReLU layer to output exactly 0. If the candidate activation mechanism is able to learn to simply output zeros when the appropriate input is zero, then if the gate is able to stay constant the task is essentially solved – there should be at least one hidden unit that will contain a value proportional to the correct answer. Therefore, we would expect all gated units with an additive state update to perform excellently on this task, those with a rectified state update only more so.

The performance of the LSTM is somewhat disappointing. Of all the models it has the most parameters and should have the most expressive power. While others have shown it is certainly capable of solving the task, the extra complexity seems to make optimising the model much more difficult. This adds weight to the hypothesis posited in section 4.2.2 that the co-operation between components that the LSTM requires to succeed is tricky to learn by gradient descent and it is preferable to avoid it if possible.

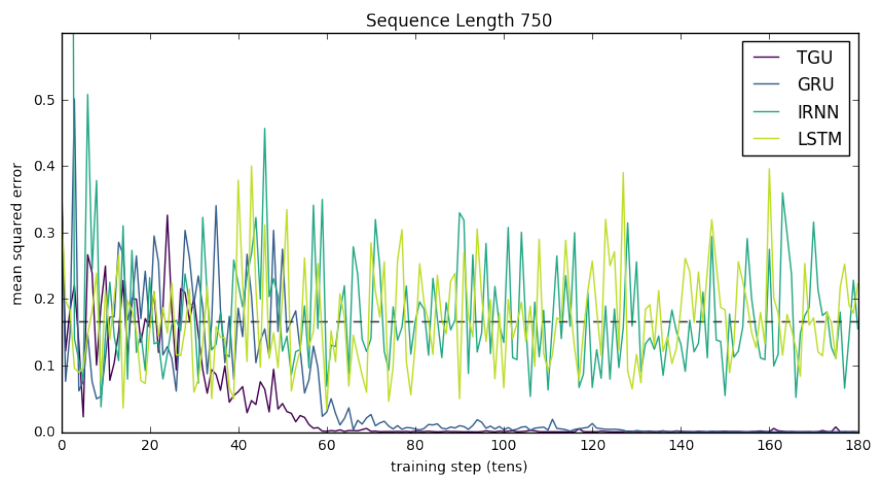
Although the convex gate is excellent at picking out individual time steps to remember, it still has the character of an exponential decay. We would then expect there might be a limit to just how far back the gradients will back-propagate and hence a limit on the longest time



(a)



(b)



(c)

Figure 5.1: Results for addition task, standard sequence lengths.

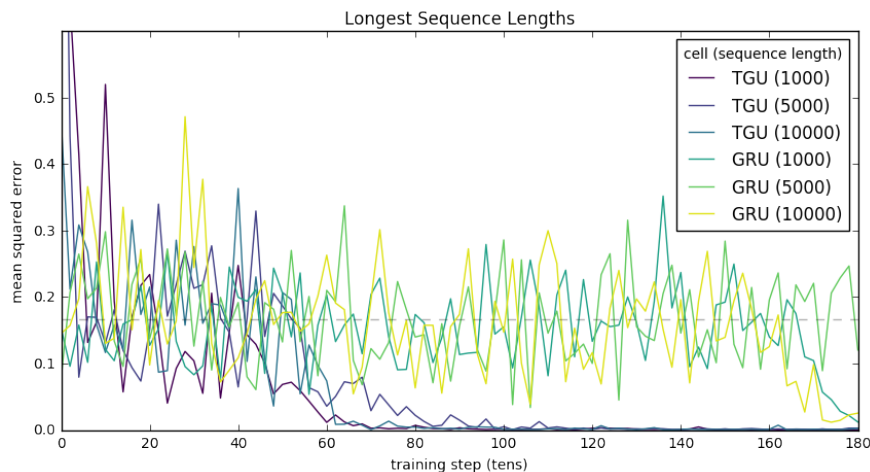


Figure 5.2: Results for addition task, very long sequences.

dependencies we can represent. Ultimately, we were unable to find a sequence long enough that the TGU was not capable of solving. Results for the TGU and GRU on extraordinarily long sequences are shown in figure 5.2. The GRU begins to beat the baseline on most of the sequences inside 2,000 updates, which is remarkable. What is more remarkable is that the TGU seems almost completely unaffected by the extreme length of the sequences, still achieving a reasonable solution inside 1,000 steps.

For the sequences with length 10,000 the expected time dependency (the number of steps between the two inputs the network has to remember) is 5,000. The network then has to preserve this value in its memory for an average of 2,500 steps. If we consider taking the RNN and unrolling it across all 10,000 timesteps, it is equivalent to a feed-forward network of extraordinary depth and shared weights at each layer.¹ In order to train it successfully we must be able to back-propagate gradients through 7,500 steps on average. While the shared weights undoubtedly assist greatly as they will allow it to generalise useful behaviours found while training on the easier sequences, it is worth observing that this is still more than 5 times the depth of some of the deepest feed-forward networks reported [46].

5.1.2 Variable Binding

This task is designed to test the architectures’ ability to store arbitrary patterns in memory and associate them with a label. We term this “variable binding” as to solve the task the network needs to be capable of associating an arbitrary value with a specific label. This task is a natural sub-problem of many large scale problem domains to which RNNs have already been successfully applied. An example is translation – in an end-to-end translation system, the network must perform tasks such as detecting subjects and objects of sentences, store them appropriately and reproduce them at an arbitrary time. **can we cite gmarcus and say that humans tend to prefer identity mappings, so maybe our nets should (when trying to do things humans are good at eg language)?**

Curiously, we could not find a satisfactory task in the literature which exercised this ability. Close alternatives include the copy task in [44] and the variants considered in [30], but these are concerned with recalling the order in which symbols arrive. The “variable assignment” task used to evaluate the Associative LSTM [19] has the same aim in mind, but

¹As the error is only computed on the final output for this task, the analogy is very close – the RNN is precisely a deep feed-forward network with shared weights and additional inputs at each layer.

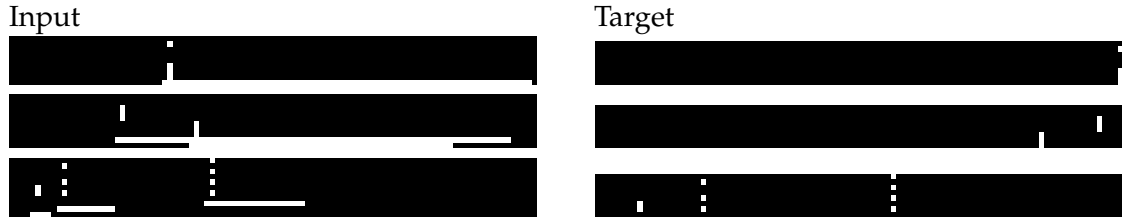


Figure 5.3: Example input/target pairs for the three different variable binding tasks. All have length 100, with 8 bit patterns – all that differs is the number of patterns the network may have to remember at once.



Figure 5.4: Outputs from TGU models that that achieved good results.

it is more complex as the assignment and retrieval instructions are spread over multiple time steps.

Task

We propose the following synthetic task to exercise the ability of the network to store and retrieve arbitrary patterns: inputs and targets are sequences of length T . Inputs at each time step are binary vectors of size $D + N$. The first D elements of the input are the “pattern bits” while the remaining N are the “label bits”. N different binary patterns are presented to the network at randomly chosen time steps in the first half of the sequence, the pattern bits are zero at all other times. Immediately before a pattern is presented, one of the label bits is set. The label bit is then held until a randomly chosen time step where it is cleared. At this point, the network must output the corresponding pattern. Label bits are never reused. The target sequence to match consists of vectors of size D , containing the appropriate pattern at the moment each of the label bits switches off and zeros elsewhere. Figure 5.3 shows example input/target pairs for various N (best viewed without antialiasing), while algorithm 2 in the appendix demonstrates how to produce appropriate data.

This task involves a number of interesting concepts. A correct solution will need to learn a transformation from input patterns to hidden states and then learn the inverse transform to

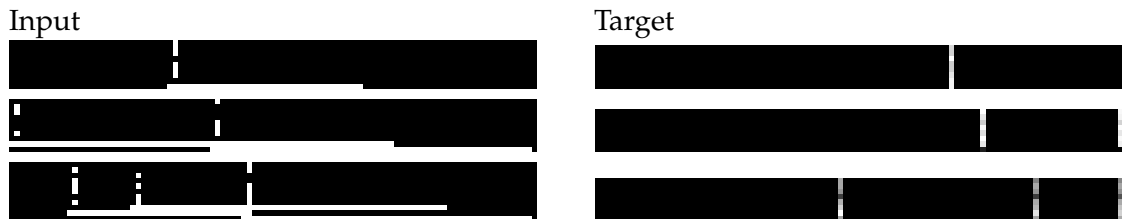


Figure 5.5: Output from GRUs demonstrating the baseline behaviour – the timing is correct but the output is not related to the input.

produce the correct output. The requirement of invertibility suggests a strong reason for the success of the TGU with linear candidate updates, as the correct output can be reconstructed by applying the inverse affine transformation. **we ensure overspecified, because learning identities is hard?**

Storing and recalling multiple distinct items also requires the ability to be selective about updating the hidden state. **gates are good for this**

Experiment Setup

For this task we found the same hyper parameters to be the best for all models (that is to say, grid search was unable to find any with demonstrably better performance). Models were trained with a batch size of 32 using ADAM with a base learning rate of 0.01. We ceased training after 5000 steps – although some models were only just beginning to learn by this time, it is sufficient to see the difference in how easily they are able to arrive at a solution.

As the targets are binary vectors we use a sigmoid output layer and train to minimise the sum of the cross entropy at each output at each time step. For all results reported we use eight bit patterns and test on sequences of length 100, 200 and 500 with one, two and three patterns leading to nine different experiments.

The baseline behaviour which all architectures (including the Vanilla RNN) achieved rapidly is to wait until any of the label bits turns off and output a pattern with mean 0.5. Guessing 0.5 (or randomly with probability 0.5) achieves a cross entropy of $-\log 0.5 = 0.69315$, so this will achieve a total loss on sequences with N eight bit patterns of $-8N \log 0.5$. Examples of a networks exhibiting this behaviour can be found in table 5.5.

As the aim is to assess the ability of the networks to solve the task, we train in an online fashion. This means that rather than generate a training and test set, we generate new sequences as required and only report the training error as in this case it is equivalent to generalisation performance.

We set the number of hidden units to ten times the number patterns to remember for each architecture. This leads to considerable difference in the number of parameters, but this task is primarily concerned with the ability of an RNN to use its memory sensibly, so it is fairest to ensure they all have the same amount of memory. As the patterns are all eight bits, have ten hidden units per pattern ensures it is possible for the networks to store the inputs directly and still have a couple of units to track the state of the label bits. As the patterns are binary, it is very easy for the networks to also learn a compressed representation, but we decided not to force them into this behaviour to keep the task focused.

Results

We report results for the LSTM, GRU, IRNN and TGU. For the TGU we incorporate the biases into the decomposition (we denote this TGU-C) and use linear activations. This was found to give the most consistent results. We also tested the Vanilla RNN, which was unable to beat the baseline in any tests and several TGU variants – for clarity we only present the results of architectures which showed at least some promise. We were unable to reliably train IRNNs on sequences longer than 100, despite heavy gradient clipping the gradients consistently exploding leading to numerical instability. Due to this we only report IRNN results on the tasks with sequence length 100.

While results were consistent, the training curves are very noisy – this is a very difficult task and none of the models are completely immune to the various problems that plague RNNs in general. To obtain a clearer picture, we plot the median of five runs.

This task becomes harder both with increased length and more patterns to remember. Increasing the length increases the length of time patterns must be remembered without degrading as well as increasing the likelihood of vanishing gradients interfering with training. Increasing the number of patterns to remember means that the RNN must learn to utilise its states in significantly more complex ways.

With only a single pattern to remember, the TGU dramatically outperformed the other architectures. The LSTM was the most consistent of the other architectures tested, although in every test by the time it had begun to escape from the baseline the TGU had already converged. The GRU was able to make surprisingly little progress on this task.

Increasing the number of items to remember, even just to two, appears to make the task considerably more difficult. This perhaps highlights an issue with the manner in which RNNs interact with their memory – they are forced to attend to it all at once. If we consider a TGU-style RNN in which the candidate state does not depend on the previous hidden states, it must learn a redundant mapping to embed the input into the state-space. This is because it can not rely on the label bits to determine which of the two inputs it is processing (this would require access to the hidden states to determine which has switched on recently). If we assume different sets of hidden states are used to represent inputs with different associated labels, then the candidate update must propose to update both of them with the new input and it is up to the gate to select which update to apply. This seems like a perfectly reasonable solution, although learning such a redundant mapping is going to be considerably more challenging than simply learning the single reversible transformation required to solve the task when only a single pattern is to be remembered.

While LSTMs and GRUs are capable of this behaviour as well as more complex solutions, the TGU is consistently able to outperform them, except on the hardest of the tasks. We suggest that restricting the number of possible solution mechanisms in this way helps the TGU find such solutions.

With two patterns to remember the TGU still makes the most progress below the baseline in the time allotted. On the longer sequences progress is limited and none of the architectures are close to converging although the TGU is the only architecture to make notable progress.

When the number of patterns is increased to three all the architectures struggle. With a sequence length of 100, there is some separation – the TGU makes the most progress although the GRU and LSTM also begin to progress beyond the baseline. On the longer sequences, there is less of a clear result and when the sequence length is up to 500, no models were capable of surpassing the baseline even when allowed to train for significantly longer than reported in figure 5.8.

These results affirm that the TGU is competitive with the state of the art approaches. They also suggest that the fundamental change – focusing the expressive power in the gate rather than the candidate update – is not misguided and makes a fundamental difference to the manner in which the network attends to its state.

5.1.3 Sequential MNIST

edit to reflect we only care about pmnist? A pair of recently proposed tasks which have gathered significant popularity lately are based around attempting to classify the MNIST handwritten image dataset [57]. To turn two dimensional images into a sequence task they are flattened in one of two ways: either in scan-line order beginning from the top-left going one row at a time or simply by taking a random permutation [56].

This gives two tasks in which the input is a one dimensional time series to be classified. The MNIST digits are 28 by 28 pixels, so the input sequences are 784 steps long. As the majority of the pixels are black with the activity focused in the central region, the majority

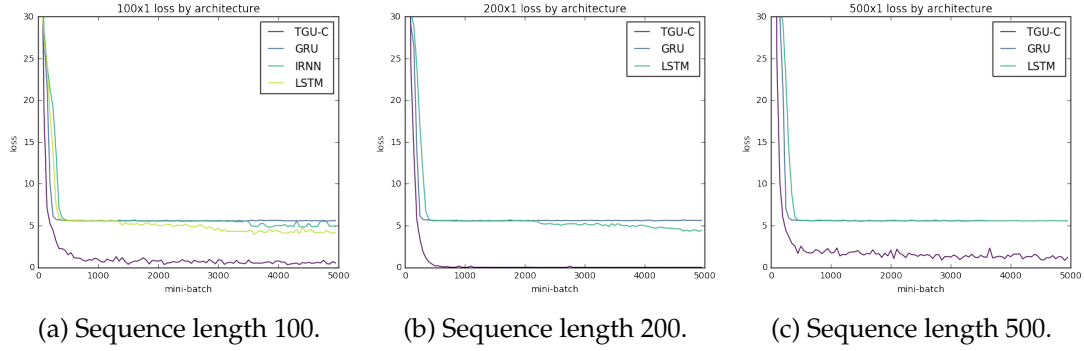


Figure 5.6: Variable binding results for sequences containing 1 pattern to remember. The baseline loss for this task is 5.5452.

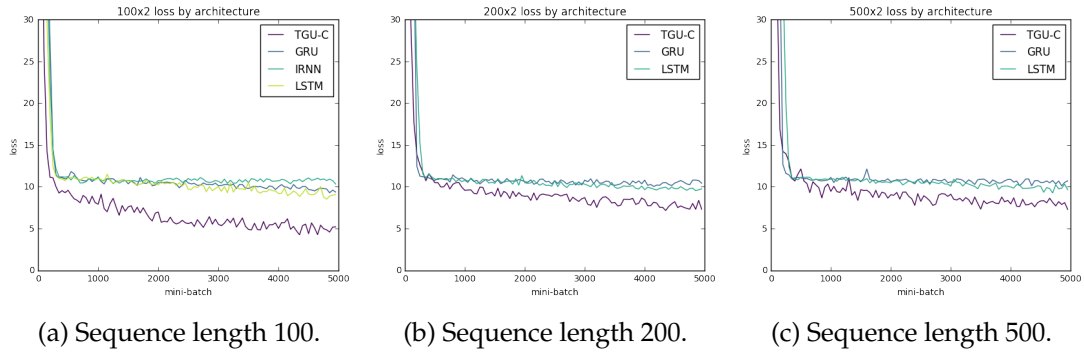


Figure 5.7: Variable binding results for sequences containing 2 patterns to remember. The baseline loss for this task is 11.0904.

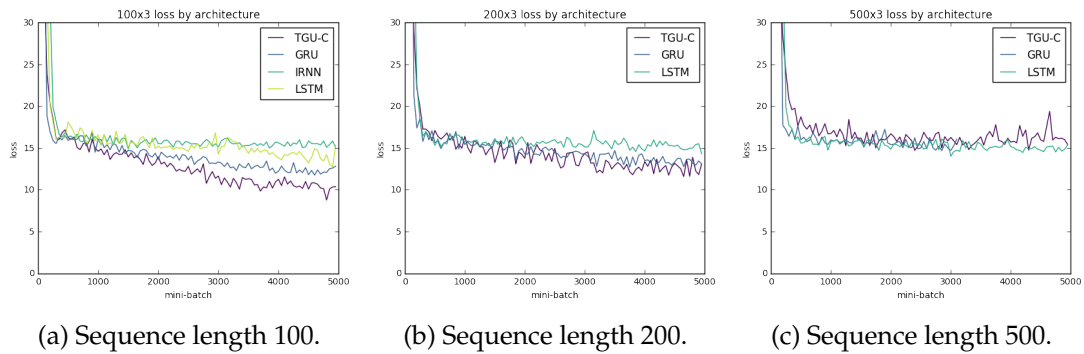


Figure 5.8: Variable binding results for sequences containing 3 patterns to remember. The baseline loss for this task is 16.6355.

of the information in the scan-line version is focused in the middle. Contrastingly, taking a random permutation of the indices is highly likely to induce very long time dependencies, in addition to removing any temporal correlations in the sequence by destroying the morphology. Consequently the permuted version is significantly more challenging.

Although this has been widely used recently to benchmark architectural advances in RNNs (see, for example [107, 5, 26, 71, 18] solely from 2016), we argue that it does not necessarily provide a useful test of a model’s capabilities. We show this by achieving excellent performance on the more challenging permuted task with a models which fail miserably in real world applications.

This is achieved by exploiting the fact that although the permutation introduces artificially long time dependencies the underlying data is very simple – a simple k-Nearest Neighbour classifier achieves 95% accuracy [57]. Combining the simplicity of the class boundaries with the sparsity of the representation in pixel-space and we can see that it should be possible to succeed simply by accumulating information at each time step, building a compressed representation which is then classified at the final time step.

Proposed Architecture

We propose an architecture specifically for this task which is a variant on GMR proposed earlier, modified to allow for direct accumulation. Further, we enforce this accumulative behaviour by applying a linear rectifier to the state update. The hidden states are computed by:

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \rho \left(\mathbf{x}_t^\top \mathbf{W} \mathbf{h}_{t-1} + \mathbf{U} \mathbf{x}_t + \mathbf{b} \right).$$

This architecture should fail – the gradients are almost guaranteed to explode during back-propagation by the argument in section 4.2.1, it can never remove information from its states (the state updates must be non-negative). The only way it can make a context-dependent decision to ignore certain inputs is if the tensor product is able to outweigh the rest of the calculation before the non-linearity and force it to be negative. As we implement the tensor product using a tensor in the CP decomposition, we refer to this architecture as the “CP+”.

Despite its clear shortcomings, this model achieves remarkable success on the sequential MNIST tasks. As discussed above we believe this is due to shortcomings with the task rather than special properties of the above model, as the above model is clearly very limited.

A model which seems slightly less absurd allows *subtractive* forgetting:

$$\begin{aligned} \mathbf{a}_t &= \rho \left(\tilde{\mathbf{x}}_t^\top \tilde{\mathbf{W}}_a \tilde{\mathbf{y}} \right) \\ \mathbf{b}_t &= \rho \left(\tilde{\mathbf{x}}_t^\top \tilde{\mathbf{W}}_b \tilde{\mathbf{y}} \right) \\ \mathbf{h}_t &= \mathbf{h}_{t-1} + \mathbf{a}_t - \mathbf{b}_t. \end{aligned}$$

We denote this architecture “CP- Δ ” and consider only the variant in which the bias matrices are incorporated into the tensor decomposition to keep the number of parameters comparable. It appears more sensible than the CP+ as it at least allows for the possibility of a negative state update. However, this model will suffer from the same gradient issues as the CP+, although they might be expected to be less pronounced as the negative term may mollify the back-propagated error. In practice, we find this is not the case and the model requires at least as much guidance as the CP+.

To enable training of this model we re-parameterise the weight matrices with a form of *weight normalisation* [83]. Typical weight normalisation learns unconstrained weights but applies a differentiable normalisation transformation before use. Salimans and Kingma

Architecture	Test Accuracy (%)
TGU	88.2
CP+	84.7
CP- Δ	91.5
IRNN (ours)	84.0
LSTM (ours)	86.3
Vanilla (ours)	79.3
i RNN [56]	82.0
u RNN [3]	91.4
s TANH-RNN [107]	94.0
BN-LSTM [18]	95.4

Table 5.1: Test accuracy for permuted sequential MNIST.

normalise each row of the weight matrices to have an l_2 norm of 1, we achieved best results by dividing each matrix by its Frobenius norm.² Back-propagation proceeds as usual through the transformation.

Experiments

We are primarily concerned with the permuted task, although we report some on the less challenging ordering as well to indicate that our models are performing near the state-of-the-art on these tasks. Learning rate and rank were found using grid-search. We train using ADAM with a batch size of 100. for nearly all models, severe gradient clipping [78] is required to train successfully, the best value was also found using grid search.

Following [56] all models tested have a single hidden layer of 100 units. We use a grid search to find the best rank for the tensor models.

Results

Final test accuracy for the best models on the permuted task is reported in table 5.1. On the simpler version all of the interesting models achieve between 97 and 99 per cent, so we do not report further results as there are no insights to be gained. Table 5.1 also includes recent results, all of which were state-of-the-art for couple of months. Our models sit in the middle of this range of results. This is remarkable as the best performing architecture, the CP- Δ has fundamental design flaws which cause it to fail on any other dataset we attempted to train it on.

Figure 5.9 shows how the performance scales with rank of the tensor decomposition. As the CP- Δ incorporates the bias matrices into the decomposition, reducing the rank should greatly affect performance. Indeed, with a rank 1 decomposition the model is incapable of learning anything. What is more curious is that the best performing model is not the highest rank. These models are not enormous, but the largest still have up to 40,000 parameters so it is possible that the lower rank provides a regularising effect which helps the model learn a more general solution.

The absurdity of the CP+ unit is clear from figure 5.9b – it trains in an incredibly volatile fashion. Interestingly for ranks 1 and 50 the best parameters found by grid search included a gradient threshold of 10, ten times the threshold needed for the other models. As rank 50 achieved the best performance, this appears to be an example of the benefits of operating

² For a $n \times m$ matrix \mathbf{U} , the Frobenius norm is $\|\mathbf{U}\|_F = \sqrt{\sum_i^n \sum_j^m U_{ij}^2} = \|\text{vec}(\mathbf{U})\|_2$.

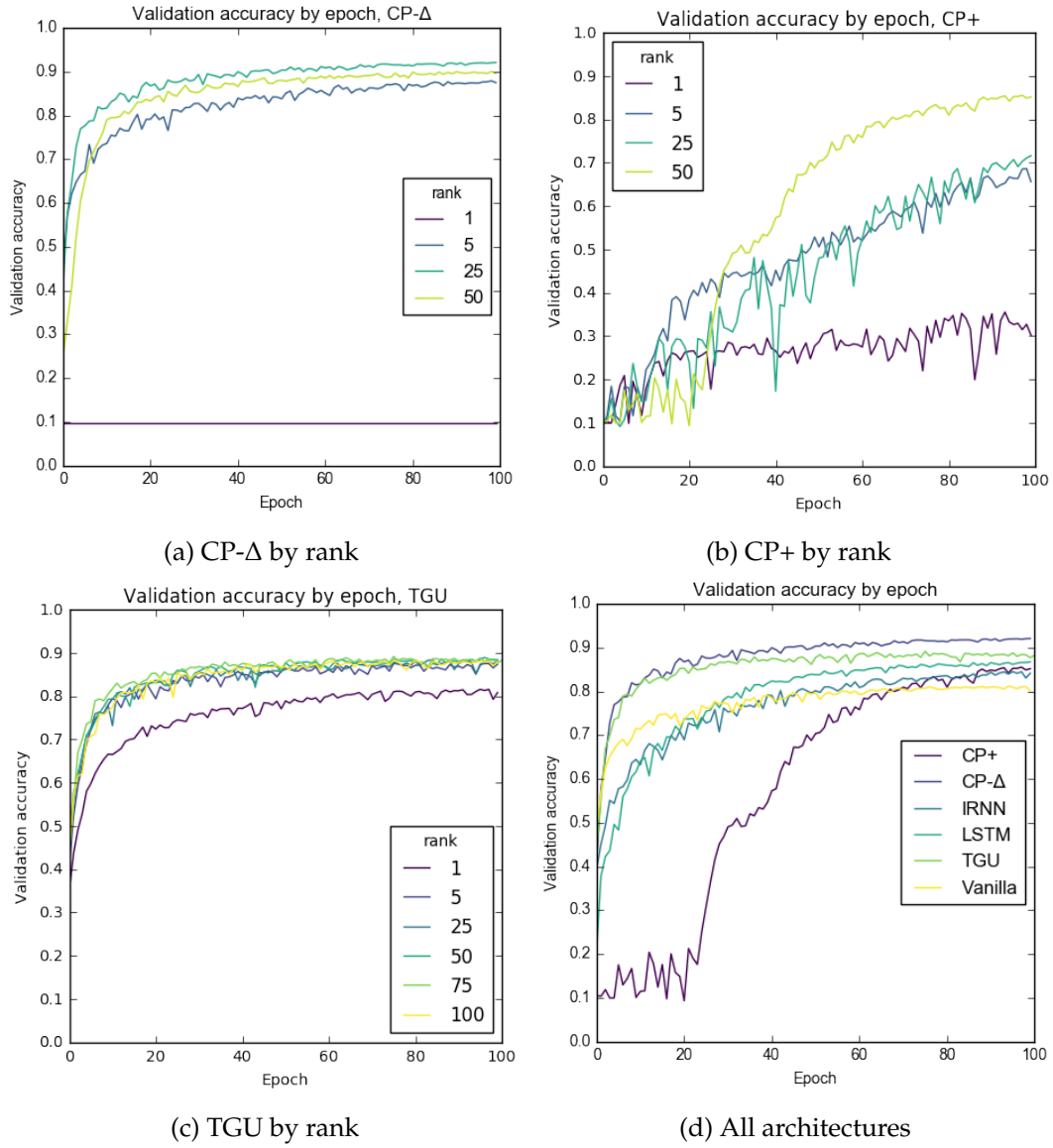


Figure 5.9: Performance (classification accuracy, %) on the validation set of permuted MNIST by rank of tensor decomposition.

at the “edge of chaos”, on the edge between vanishing and exploding gradients and states [7]. Further, that the rank 1 model learns at all is remarkable. The tensor product is this unit is severely restricted, the vast majority of its learning power resides in a feed-forward connection. This goes some way to affirming our suspicions about this task – although it exhibits long time dependencies, the benefit of using detailed contextual information at each time step is limited.

5.2 Real-world Data

We now move on to more realistic datasets. All of the following tasks are *next-step prediction* tasks – the goal of the network is to estimate the probability of the next item in the sequence given the sequence so far:

$$p(\mathbf{x}_{t+1} | \mathbf{x}_1, \dots, \mathbf{x}_t).$$

The output layer of the RNN then has the task of converting the hidden states into probabilities and so requires an appropriate non-linearity. Training proceeds by maximising the log-likelihood of the data producing the correct distribution.

Essentially this means that the target sequence is the input sequence shifted across by one. If we prepend a special *start-of-sequence* symbol to the data sequence and append an *end-of-sequence* symbol, then the inputs to the network are simply all but the *end-of-sequence* and the outputs are all but the *start-of-sequence*.

The end result is a statistical model of the sequence. We can sample from it by simply feeding it the *start-of-sequence* character and then repeatedly sampling from the distribution generated by the current hidden states and feeding the sample back in to the network. This can be an interesting way to produce novel sequences [29].

In all cases the data is split into three sets: training, test and validation. The size of the splits vary as we attempt to remain consistent with prior work. The validation set is used to monitor generalisation performance and to stop training when the model begins to overfit to the training data. The manner in which this is realised is by checking the average negative log-likelihood the model achieves after each epoch of training. If the validation performance has improved, we save a checkpoint containing the model parameters. After having trained for a fixed number of epochs the best saved model is loaded and its performance is evaluated on the test set.

The aim of this section is not to achieve state-of-the-art performance on large datasets using enormous networks. Rather we seek to ascertain whether the theoretical reasoning and intuitions developed in the preceding chapters hold in a more realistic setting. **some nice way to say this?? but the idea is that we don't want to compete with hacks but to see if the ideas and guiding principles we came up with are indeed sensible and to see whether our RNNs are competitive in a restricted setting on smaller data without fancy engineering tricks.**

5.2.1 Polyphonic Music

This task, introduced by [8] consists of modelling four datasets of polyphonic music from a score. This is potentially a challenging task – not only do the melodies and chord progression exhibit long time dependencies, the networks must learn the rules of harmony to produce appropriately consonant arrangements. The likelihood of a given note at a given time step is affected by both simultaneous information (which other notes are one at the same time) and long term information (such as which key the piece is in). Modelling accurately both these dependencies is challenging but reflects well the difficulties inherent in much naturally generated data.

Task

The four datasets are:

Pianomidi is a collection of classical piano scores originally sourced from <http://piano-midi.de>. We split the data according to [80]. This data uses the full range of the piano, there are 88 distinct notes present although the polyphony is limited as the pieces must all be played by a human.

Nottingham is a set of folk tunes converted to MIDI by [8],³ we use their split of the data. These tunes have the simplest structure and are generally quite short although they span a range of 63 distinct notes.

Muse comes originally from <http://www.musedata.org> and contains both orchestral and piano classical music. This dataset has a range of 84 distinct notes and due to the orchestral nature of many of the scores the polyphony can be quite high. Again we split the data according to [8].

JSB consists of all 382 chorales harmonised by J. S. Bach. This dataset has some of the lowest polyphony, at most four notes at once, although it is very consistent by its nature. This dataset also has the smallest range with a final range of 54 notes. We use the split of the data given in [2].

The goal is to maximise the likelihood the model assigns to the data. This is done by minimising the negative log-likelihood. In general, Nottingham is the simplest of the datasets primarily due to the simple harmonic structure.

Experiments

The data is converted from MIDI files (which contain note on/off events) to an appropriate format by sampling the notes present at regular time intervals. This gives binary vectors in which each element indicates with a one or a zero whether a particular note is active or not.⁴ For each dataset, we fix the size of these vectors to the range of notes present in the dataset (after transposing to C major or C minor, equivalent to centering the data).

Previous work is somewhat inconsistent with the size of these vectors. Of those that report it, [8] fix the size of the vectors for all datasets to the range of 88 notes available on a piano while [15] end up with much larger vectors although it is not clear why. We fix the size of the inputs to the range of the notes present in the data. Although it could be argued it would be preferable and more flexible to use a more generous range, when the data is music written for specific instruments it only makes sense to train the model to attend to an appropriate range for those instruments. In general this means results can not be directly compared although it is worth noting they are similar, with nearly all of our RNNs surpassing [15] and achieving performance similar to the RNNs evaluated in [8]. Further, the more unexpected elements of the results such as the performance of the Vanilla RNN on Nottingham are present in others' results. Given the differences in setup, we would expect our RNNs to fare slightly worse as using larger data vectors means many of the output units are redundant and never used. The network should learn this quite quickly which will bring reduce the average loss.

All networks used a single hidden layer. The number of hidden nodes and the rank of the tensor decomposition was chosen so that each network would have as close to 20,000

³ The originals are available at: <http://ifdo.ca/~seymour/nottingham/nottingham.html>

⁴ Appropriately pre-processed data can be downloaded from <http://www-etud.iro.umontreal.ca/~boulanni/icml2012> courtesy of Boulanger-Lewandowski et al. [8].

Model	hidden units/rank (parameters)			
Vanilla	96 (19927)			
LSTM	44 (20075)			
GRU	52 (19763)			
GMR	95/1 (19870)	71/35 (19872)	58/58 (19775)	45/90 (20125)
GMR-C	349/1 (20005)	106/52 (20142)	76/76 (20119)	53/106 (20248)
TGU	80/1 (20030)	63/31 (20156)	53/53 (20248)	41/82 (19817)
TGU-C	176/1 (20000)	88/44 (20075)	66/66 (19855)	48/96 (20071)

Table 5.2: Size of models for polyphonic music modelling. Architectures with -C appended have the bias matrices combined with the decomposition. Parameters are reported for inputs of size 54 as per the JSB dataset. Rank is only reported if applicable.

	Nottingham		JSB		Pianomidi		Muse	
	<i>Train</i>	<i>Test</i>	<i>Train</i>	<i>Test</i>	<i>Train</i>	<i>Test</i>	<i>Train</i>	<i>Test</i>
Vanilla	3.3278	3.8871	8.0284	8.6381	7.3874	7.8093	6.8954	7.3619
LSTM	3.5826	4.0124	8.2058	8.6633	7.3730	7.7542	7.0777	7.4159
GRU	3.5876	4.0157	8.1689	8.6027	7.3039	7.7652	7.0001	7.3824
GMR	3.3940	3.9055	8.0526	8.6084	7.3598	7.7676	6.8864	7.3089
GMR-C	3.7902	4.1350	8.0093	8.5369	7.4244	7.8615	6.9135	7.4296
TGU	3.9652	4.2140	7.9357	8.6060	7.2132	7.7161	7.4524	7.6879
TGU-C	4.0035	4.2531	7.9731	8.5307	7.2748	7.7003	7.4551	7.6719

Table 5.3: Results on polyphonic music datasets. Numbers are average negative log-likelihood, lower is better. “-C” appended to the tensor units indicates the bias matrices are folded into the decomposition, otherwise they are separate.

parameters when applied to the JSB dataset as possible. This was done to help make the comparison more fair so that the comparison is not skewed by models such as the LSTM which have a large number of parameters per hidden unit. The rank was set as either 1 or a multiple of the hidden states, chosen from $\{1/2, 1, 2\}$. Table 5.2 shows the number of hidden units and the rank of the tensor decomposition for the models evaluated. The networks used a sigmoid output layer to appropriately model the fact that the multiple output units may need to be on at once.

A preliminary grid search was undertaken on a subset of the architectures to determine sensible ranges for the final experiments. Models reported were trained using ADAM with a grid search over the following hyper-parameters: base learning rate from $\{0.01, 0.001\}$, maximum length of back-propagation through time chosen from $\{75, 100\}$ steps and rank chosen from the options in table 5.2. Note that to keep the number of parameters constant, changing the rank necessitates changing the number of hidden units as well. The batch size was set to 8. For nearly all architectures best performance was achieved with the lower learning rate while the best maximum length was task-dependent.

Results

Results for the four datasets are found in table 5.3. We report the average negative log-likelihood the best early-stopped models assign the data vectors over the training, test and validation sets.

In general the tensor units performed better than their counterparts except on Nottingham, the simplest of the datasets. This dataset consists of folk songs which are often very formulaic

in their harmony. Correspondingly, simple models seem to be able to perform very well. It also seems to be the case that there is something about the structure of the data which suits the Vanilla RNN – the best GMR (which was the closest to the Vanilla RNN out of all architectures) had a rank 1 tensor decomposition. This means that it was functionally nearly identical to a Vanilla RNN.

In fact, all of the best GMR models used a rank 1 tensor decomposition. This is interesting as it suggests the extra complication of the tensor product was only hindered learning on this data. An alternative conclusion that may be drawn is that it is beneficial to have more hidden units, even if that comes at the cost of decreased representational power. This effect was only notable in the GMR, so it is a distinct possibility given its general inability to modify its state in the complex manner captured by the gated architectures.

still more to say about rank

5.2.2 PTB

The next experiment aims to test whether reducing the number of parameters in the model by lowering the rank of the tensor decomposition can have a useful regularising effect. For this we use the Penn Treebank [61] and attempt to model it by word using the next-step prediction framework outlined earlier. This is a relatively small dataset and overfitting becomes a serious challenge. Due to this the dataset has been popular for testing regularisation approaches for recurrent neural networks [106, 25].

We are concerned with how the tensor that computes the gate values in the TGU affects the ability of the network to make use of its memory. To test this we fix the number of hidden units and vary the rank of the tensor.

Experiment Setup

The data is pre-processed following [106] by replacing all but the most common 10,000 words with an ‘unknown’ symbol. We train again using ADAM with a batch size of 20. Base learning rate for each model was found via grid search.

All networks trained had 128 hidden units – for this experiment we do not attempt to keep the number of trainable parameters equal between models. We test four variations of the TGU: with and without combined biases with with linear or rectified linear candidate updates. For comparison we also train an LSTM, GRU and Vanilla RNN. We train for ten epochs and save the model from the best iteration according to validation performance to be evaluated against the test data.

Results

For this task we report the average word-level perplexity over the test set. The perplexity is the exponentiation of the cross entropy between the predicted distribution and the target distribution (lower is better).

Overall results are presented in table 5.4. The best performance across all models was achieved by the TGU with linear candidate updates and separate biases. Three out of four of the TGUs surpassed the LSTM with linear candidate updates having a clear advantage.

Figure 5.10 shows results of the TGUs organised by rank. For all models there is a clear benefit to constraining the tensor to half the number of hidden states. When the biases are combined the difference between rank 64 and 128 is less clear cut, which is to be expected as the expressive power of the tensor unit is much more constrained by the rank. That none of the models with separate biases improve performance with a rank below 64 indicates that

Architecture (rank)	<i>Train</i>	<i>Valid</i>	<i>Test</i>
TGU (64)	82.86	139.54	131.42
TGU-C (64)	92.73	147.90	139.15
lin-TGU (64)	91.81	127.80	121.93
lin-TGU-C (128)	99.74	133.77	127.56
LSTM	78.74	144.75	134.79
GRU	79.89	133.55	127.56
Vanilla	85.10	164.81	156.63

Table 5.4: Results of best early-stopped models on the Penn Treebank

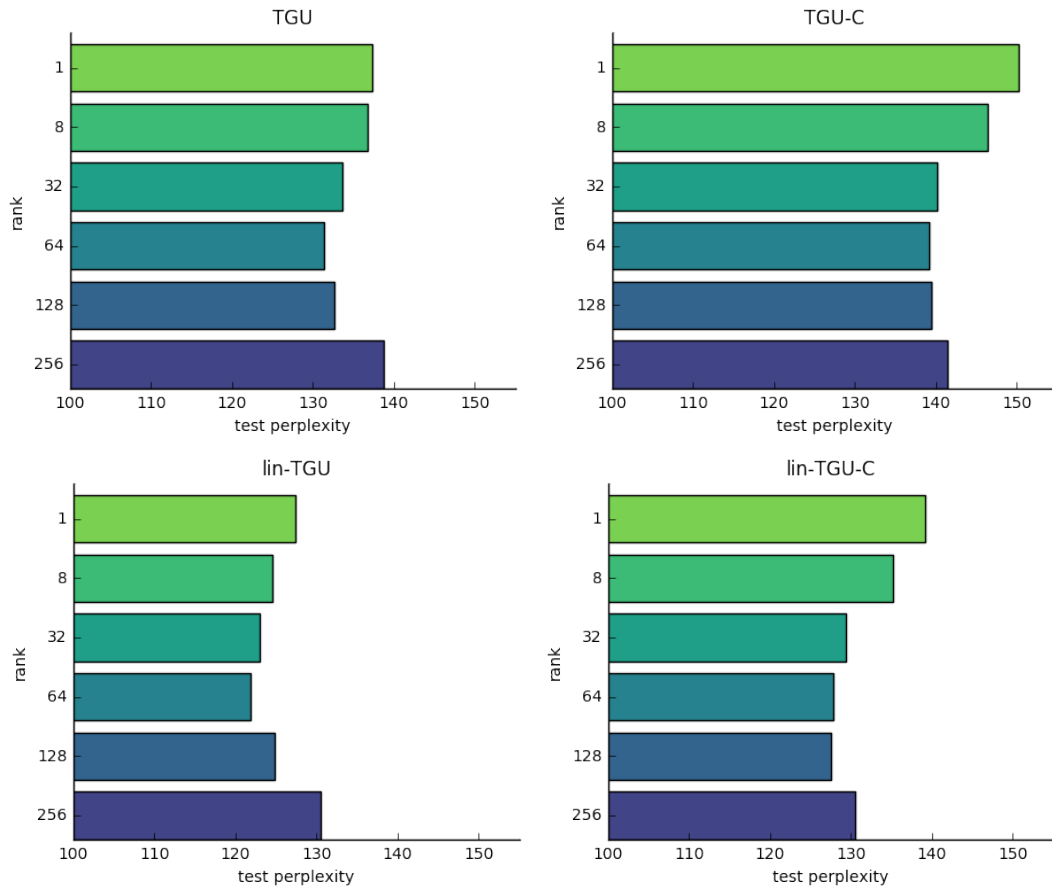


Figure 5.10: Test perplexity of best early stopped models for each type of TGU for a range of ranks. Networks had 128 hidden units.

the tensor is performing a valuable role as without the tensor the gate unit of these models reduces to the standard LSTM/GRU gating unit.

5.2.3 War and Peace

We run one final experiment to test whether the regularising effect observed so far are due to limiting the expressive power of the tensor or simply reducing the parameters. We do this by adjusting the number of hidden units so that all models have the same number of trainable parameters regardless of rank.

As data for this test we use the novel War and Peace by Leo Tolstoy to build character level language models. This is mostly in English and in order to succeed the models must learn to model long term temporal dependencies such as opening and closing quotations marks, correct punctuation and grammar as well as higher level language concepts.

Experiment Setup

Experiments are set up very similarly to [50]. Accordingly we split the 3,258,246 characters into training/validation/test sets with an 80/10/10 % ratio, giving approximately 300,000 characters in the training and validation sets. The book contained a vocabulary of 83 characters. Again we train using ADAM, we found a learning rate of 0.001 to be robust for all models. We use a batch size of 100 and limit back-propagation of gradients to 100 steps. Models were trained for 50 epochs and the test error was calculated using the best model based on validation error as for the other datasets.

We found even the smallest models overfit dramatically even after one or two epochs. To even approach reasonable baseline performance, we introduce a bottleneck before the network in the form of very small character embeddings. This is equivalent to encoding the input symbols as one-hot vectors (size 83×1) and multiplying by a 8×83 matrix to embed the symbols into an 8 dimensional space. This matrix is learned jointly with the model. To further hold off overfitting we apply dropout [88] to the embeddings before feeding them to the network.

All models were constrained to have a budget of 25,000 parameters (including the output layer but not including the embedding matrix as this is fixed size). This is quite small, it corresponds to a LSTM with 65 hidden units or a Vanilla RNN with 118 on this data. The rank one GMR with combined biases, by contrast, was able to have 293 hidden units, although it proved incapable of making use of the them. Vanilla RNNs perform surprisingly well on this task, so we test the GMR in addition to the TGU. As we are interested in the effect of low-rank tensor decompositions on expressive power, we only report results for models with combined biases. We found linear candidate activations consistently outperformed rectifiers for the TGU.

For the tensor units we choose rank proportional to the number of hidden units, with the proportion giving a measure of the expressive power.

Results

Training and validation cross entropy are plotted in figure 5.11, averaged per character. The tensor models both show remarkable resistance to overfitting compared to the other architectures. What is particularly interesting is that the validation error does not begin to rise seriously at all, regardless of the rank. Aside from rank one, which struggles to learn at all, the models benefit slightly from lower rank with best test performance being achieved by both tensor models with rank equal to a quarter of their number of hidden units.

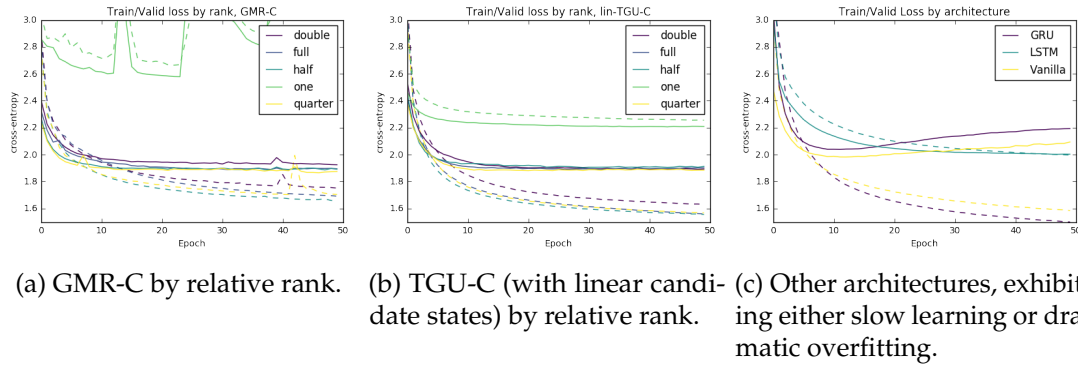


Figure 5.11: Training and validation curves on War and Peace. Dashed lines indicate training loss, solid lines are the loss on the validation set. ‘Full’ refers to rank equal to the number of hidden units, ‘one’ simply means rank 1.

Architecture	Test Cross-Entropy
GMR-C	2.018
lin-TGU-C	2.052
Vanilla	2.128
GRU	2.1621
LSTM	2.1271

Table 5.5: Test set cross entropy (per character) for the best War and Peace models.

Perhaps more unexpectedly, even the tensor model with an excessive rank of twice the number of hidden units still fails to overfit as dramatically as the LSTM or the Vanilla RNN. This seems to suggest that the tensor units work particularly well with the small embeddings and dropout which were primarily applied to allow the other architectures to report validation cross entropies below three or four at the end of training.

Chapter 6

Conclusions

The conclusions are presented in this Chapter.

Appendix A

Additional Proofs

Proposition A.0.1 (Identity Tensor). $\mathcal{H} \in \mathbb{R}^{N \times N \times N}$ such that

$$\mathbf{x}^\top \mathcal{H} \mathbf{y} = \mathbf{x} \odot \mathbf{y}, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N$$

implies

$$H_{ijk} = \begin{cases} 1 & \text{if } i = j = k \\ 0 & \text{otherwise.} \end{cases}$$

Proof. We prove briefly, by inspecting one component of the result. Let $\mathbf{z} = \mathbf{x}^\top \mathcal{H} \mathbf{y}$. Then

$$\begin{aligned} z_j &= \mathbf{x}^\top \mathbf{H}_{\cdot j} \mathbf{y} \\ &= \sum_i^N \sum_k^N x_i H_{ijk} y_k \end{aligned}$$

If $z_j = x_j y_j$ as in the elementwise product, then it is clear we want H_{ijk} to be 1 if $i = j = k$. Further, if we ensure H_{ijk} is 0 when this is not the case we can see that the rest of the terms in the sums will disappear. \square

Appendix B

Algorithms for Synthetic Tasks

B.1 Addition

Data for the addition task can be generated by algorithm 1. This algorithm produces a single item, it could also be done in batches with slight adjustments.

Algorithm 1: Generating data for addition task
--

<p>Input: Integer T Output: Problem instance $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T), y$</p> <p>Sample integer $i \sim U[1, (T/2) - 1]$ Sample integer $j \sim U[T/2, T]$ $y \leftarrow 0$ for $t \in \{1, \dots, T\}$ do Sample float $x_{t1} \sim U[0, 1]$ if $t = i$ or $t = j$ then $y \leftarrow y + x_{t1}$ $x_{t2} \leftarrow 1$ else $x_{t2} \leftarrow 0$ return $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T), y$</p>

B.2 Variable binding

Algorithm 2 shows how to generate a single data item for this task. This algorithm is indicative only, there are ways to be more efficient especially if it has to be implemented as a computation graph, for example in Tensorflow.

Algorithm 2: Generating data for variable binding**Input:** Integers T, D, N **Output:** Problem instance $(\mathbf{x}_1, \dots, \mathbf{x}_T), (\mathbf{y}_1, \dots, \mathbf{y}_T)$ initialise $(\mathbf{x}_1, \dots, \mathbf{x}_T), (\mathbf{y}_1, \dots, \mathbf{y}_T)$ to zeros**for** $i \in \{1, \dots, N\}$ **do** Sample integer $j \sim U[1, (T/2) - 1]$

/* start position */

 Sample integer $k \sim U[j + 1, T]$

/* end position */

 $\mathbf{z} \leftarrow$ random D bit binary pattern $\mathbf{y}_{k+1} \leftarrow \mathbf{z}$ Set first D bits of \mathbf{x}_{i+1} to \mathbf{z} **for** $l \in \{j, \dots, k\}$ **do** set label bits $x_{l,D+i} \leftarrow 1$ **return** $(\mathbf{x}_1, \dots, \mathbf{x}_T), (\mathbf{y}_1, \dots, \mathbf{y}_T)$

Bibliography

- [1] ABADI, M., ET AL. TensorFlow: Large-scale machine learning on heterogeneous systems. *arXiv preprint* (2015). arXiv: arXiv:1603.04467v2.
- [2] ALLAN, M., AND WILLIAMS, C. K. I. Harmonising Chorales by Probabilistic Inference. In: *Advances in Neural Information Processing Systems*, 17. 2005, 25–32.
- [3] ARJOVSKY, M., SHAH, A., AND BENGIO, Y. Unitary Evolution Recurrent Neural Networks. *arXiv* (Nov. 2015), 1–11. arXiv: 1511.06464.
- [4] BALDUZZI, D., AND GHIFARY, M. Strongly-Typed Recurrent Neural Networks. In: *International Conference on Machine Learning - ICML 2014*. 2016, 9. arXiv: 1602.02218.
- [5] BARONE, A. V. M. Low-rank passthrough neural networks. *Arxiv* (Mar. 2016), 16. arXiv: 1603.03116.
- [6] BENGIO, Y. Learning Long-Term Dependencies with Gradient Descent is Difficult. *IEEE Transactions on Neural Networks* 5, 2 (1994), 157–166.
- [7] BERTSCHINGER, N., AND NATSCHLGER, T. Real-time computation at the edge of chaos in recurrent neural networks. *Neural computation* 16, 7 (2004), 1413–36.
- [8] BOULANGER-LEWANDOWSKI, N., VINCENT, P., AND BENGIO, Y. Modeling Temporal Dependencies in High-Dimensional Sequences: Application to Polyphonic Music Generation and Transcription. *Proceedings of the 29th International Conference on Machine Learning (ICML-12) Cd* (June 2012), 1159–1166. arXiv: 1206.6392.
- [9] CARROLL, J. D., AND CHANG, J. J. Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition. *Psychometrika* 35, 3 (Sept. 1970), 283–319.
- [10] CHAN, W., ET AL. Listen, attend and spell. *arXiv preprint* (2015), 1–16. arXiv: 1508.01211.
- [11] CHO, K., ET AL. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (June 2014), 1724–1734. arXiv: 1406.1078.
- [12] CHO, K., ET AL. On the Properties of Neural Machine Translation: EncoderDecoder Approaches. *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation* (Sept. 2014), 103–111. arXiv: arXiv:1409.1259v2.
- [13] CHOI, K., ET AL. Convolutional Recurrent Neural Networks for Music Classification. *arXiv preprint* (Sept. 2016). arXiv: 1609.04243.
- [14] CHUNG, J., AHN, S., AND BENGIO, Y. Hierarchical Multiscale Recurrent Neural Networks. *arXiv preprint* (2016). arXiv: 1609.01704.
- [15] CHUNG, J., ET AL. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. *arXiv* (Dec. 2014), 1–9. arXiv: 1412.3555v1.

- [16] CHUNG, J., ET AL. Gated feedback recurrent neural networks. *Proceedings of the 32nd International Conference on Machine Learning* 37 (2015), 2067–2075. arXiv: 1502.02367.
- [17] CICHOCKI, A., ET AL. Low-Rank Tensor Networks for Dimensionality Reduction and Large-Scale Optimization Problems: Perspectives and Challenges PART 1. *arXiv preprint* (2016), 100. arXiv: 1609.00893.
- [18] COOIJMANS, T., ET AL. Recurrent Batch Normalization. *arXiv preprint* (2016), 1–10. arXiv: 1603.09025.
- [19] DANIHELKA, I., ET AL. Associative Long Short-Term Memory. *arXiv* (Feb. 2016). arXiv: 1602.03032.
- [20] DHINGRA, B., ET AL. Gated-Attention Readers for Text Comprehension. *arXiv* (June 2016). arXiv: 1606.01549.
- [21] DUVENAUD, D., ET AL. Avoiding pathologies in very deep networks. *AISTATS* (Feb. 2014), 202–210. arXiv: 1402.5836.
- [22] DZMITRY BAHDANA, ET AL. Neural Machine Translation By Jointly Learning To Align and Translate. *International Conference On Learning Representations* (2015), 1–15. arXiv: 1409.0473.
- [23] ELDAN, R., AND SHAMIR, O. The Power of Depth for Feedforward Neural Networks. In: *29th Annual Conference on Learning Theory*. Dec. 2016, 907–940. arXiv: 1512.03965.
- [24] ELMAN, J. I. Finding Structure in Time. *COGNITIVE SCIENCE* 14 (1990), 179–211.
- [25] GAL, Y. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. *arXiv preprint* (2015). arXiv: 1512.05287.
- [26] GAO, Y., AND GLOWACKA, D. Deep Gate Recurrent Neural Network. *arXiv preprint* 1 (2016), 7. arXiv: 1604.02910.
- [27] GLOROT, X., AND BENGIO, Y. Understanding the difficulty of training deep feed-forward neural networks. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)* 9 (2010), 249–256.
- [28] GOODFELLOW, I. J., ET AL. Maxout Networks. *Proceedings of the 30th International Conference on Machine Learning (ICML)* 28 (2013), 1319–1327. arXiv: 1302.4389.
- [29] GRAVES, A. Generating sequences with recurrent neural networks. *arXiv preprint* (2013), 1–43. arXiv: arXiv:1308.0850v5.
- [30] GRAVES, A., WAYNE, G., AND DANIHELKA, I. Neural Turing Machines. *arXiv preprint* (2014), 1–26. arXiv: arXiv:1410.5401v2.
- [31] GRAVES, A., ET AL. Connectionist Temporal Classification : Labelling Unsegmented Sequence Data with Recurrent Neural Networks. *Proceedings of the 23rd international conference on Machine Learning* (2006), 369–376.
- [32] GREFFENSTETTE, E., ET AL. Learning to Transduce with Unbounded Memory. *arXiv preprint* (2015), 12. arXiv: 1506.02516.
- [33] GREFF, K., ET AL. LSTM: A Search Space Odyssey. *IEEE Transactions on Neural Networks and Learning Systems* (2016). arXiv: 1503.04069.
- [34] GREGOR, K., ET AL. DRAW: A Recurrent Neural Network For Image Generation. *Icml-2015* (Feb. 2015), 1462–1471. arXiv: 1502.04623.
- [35] HARSHMAN, R. A. Foundations of the PARAFAC procedure: Models and conditions for an “explanatory” multimodal factor analysis. *UCLA Working Papers in Phonetics* 16, 10 (1970), 1–84.

- [36] HE, K., ET AL. Deep Residual Learning for Image Recognition. *arXiv* (Dec. 2015). arXiv: 1512.03385.
- [37] HE, K., ET AL. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. *arXiv preprint* (2015). arXiv: 1502.01852.
- [38] HENAFF, M., SZLAM, A., AND LECUN, Y. Orthogonal RNNs and Long-Memory Tasks. *arxiv* (2016). arXiv: 1602.06662.
- [39] HIHI, S. E., AND BENGIO, Y. Hierarchical Recurrent Neural Networks for Long-Term Dependencies. *Nips* (1995), 493–499.
- [40] HINTON, G. E. *A Parallel Computation that Assigns Canonical Object-Based Frames of Reference*. 1981.
- [41] HINTON, G. E., AND SALAKHUTDINOV, R. R. Reducing the Dimensionality of Data with Neural Networks. *Science* 313, 5786 (2006), 504–507. arXiv: 20.
- [42] HITCHCOCK, F. L. Multiple Invariants and Generalized Rank of a P-Way Matrix or Tensor. *Journal of Mathematics and Physics* 7, 1-4 (Apr. 1928), 39–79.
- [43] HITCHCOCK, F. L. The Expression of a Tensor or a Polyadic as a Sum of Products. *Journal of Mathematics and Physics* 6, 1-4 (Apr. 1927), 164–189.
- [44] HOCHREITER, S., AND SCHMIDHUBER, J. Long Short-Term Memory. *Neural Computation* 9, 8 (Nov. 1997), 1735–1780. arXiv: 1206.2944.
- [45] HORNIK, K., STINCHCOMBE, M., AND WHITE, H. Multilayer Feedforward Networks Are Universal Function Approximators. *Neural Networks* 2 (1989), 359–366.
- [46] HUANG, G., LIU, Z., AND WEINBERGER, K. Q. Densely Connected Convolutional Networks. *arXiv preprint* (2016), 1–12. arXiv: 1608.06993.
- [47] JOULIN, A., AND MIKOLOV, T. Inferring Algorithmic Patterns with Stack-Augmented Recurrent Nets. *arXiv* (2015), 1–10. arXiv: arXiv:1503.01007v4.
- [48] JOZEFOWICZ, R., ZAREMBA, W., AND SUTSKEVER, I. An Empirical Exploration of Recurrent Network Architectures. In: *ICML*. 2015.
- [49] KAISER, ., AND SUTSKEVER, I. Neural GPUs Learn Algorithms. *arXiv:1511.08228 [cs]* (2015), 1–9. arXiv: 1511.08228.
- [50] KARPATHY, A., JOHNSON, J., AND FEI-FEI, L. Visualizing and Understanding Recurrent Networks. *ICLR* (June 2016), 1–13. arXiv: arXiv:1506.02078v1.
- [51] KINGMA, D. P., AND BA, J. L. ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION (2014). arXiv: 1412.6980.
- [52] KOLDA, T. G., AND BADER, B. W. Tensor Decompositions and Applications. *SIAM Review* 51, 3 (Aug. 2009), 455–500.
- [53] KOUTNIK, J., ET AL. A Clockwork RNN. *Proceedings of The 31st International Conference on Machine Learning* 32 (2014), 1863–1871. arXiv: arXiv:1402.3511v1.
- [54] KRUEGER, D., AND MEMISEVIC, R. Regularizing RNNs by Stabilizing Activations. *International Conference On Learning Representations* (Nov. 2016), 1–8. arXiv: 1511.08400.
- [55] KURACH, K., ANDRYCHOWICZ, M., AND SUTSKEVER, I. Neural Random-Access Machines. *ICLR* (2016), 17. arXiv: 1511.06392.
- [56] LE, Q. V., JAITLEY, N., AND HINTON, G. E. A Simple Way to Initialize Recurrent Networks of Rectified Linear Units. *arXiv preprint arXiv:1504.00941* (2015), 1–9. arXiv: arXiv:1504.00941v1.

- [57] LECUN, Y., ET AL. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 11 (1998), 2278–2324.
- [58] LORENZ, E. N. Deterministic Nonperiodic Flow. *Journal of the Atmospheric Sciences* 20, 2 (Mar. 1963), 130–141.
- [59] LUONG, M.-T., ET AL. MULTI-TASK SEQUENCE TO SEQUENCE LEARNING. In: *International Conference On Learning Representations*. Nov. 2016. arXiv: arXiv:1511.06114v4.
- [60] MAGNUS, J., AND NEUDECKER, H. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. 3rd ed. Wiley, 2007.
- [61] MARCUS, M. P., SANTORINI, B., AND MARCINKIEWICZ, M. A. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics* 19, 2 (1993), 313–330.
- [62] MARTENS, J., AND SUTSKEVER, I. Learning recurrent neural networks with Hessian-free optimization. *Proceedings of the 28th International Conference on Machine Learning, ICML 2011* (2011), 1033–1040.
- [63] MEMISEVIC, R. Learning to relate images: Mapping units, complex cells and simultaneous eigenspaces. *arXiv preprint arXiv:1110.0107* (2011), 1–32. arXiv: 1110.0107.
- [64] MEMISEVIC, R., AND HINTON, G. Unsupervised Learning of Image Transformations. *Computer Vision and Pattern Recognition, IEEE Computer Society Conference on* (2007), 1–8.
- [65] MEMISEVIC, R., AND HINTON, G. E. Learning to represent spatial transformations with factored higher-order Boltzmann machines. *Neural computation* 22, 6 (2010), 1473–1492.
- [66] MIKOLOV, T. Statistical Language Models Based on Neural Networks. PhD thesis. 2012, 1–129. arXiv: 1312.3005.
- [67] MIKOLOV, T., ET AL. Learning Longer Memory in Recurrent Neural Networks. *Iclr* (Dec. 2015), 1–9. arXiv: arXiv:1412.7753v1.
- [68] MINSKY, M., AND PAPERT, S. *Perceptrons*. M.I.T. Press, 1969.
- [69] NAIR, V., AND HINTON, G. E. Rectified Linear Units Improve Restricted Boltzmann Machines. *Proceedings of the 27th International Conference on Machine Learning* 3 (2010), 807–814.
- [70] NEYSHABUR, B., SALAKHUTDINOV, R., AND SREBRO, N. Path-SGD: Path-Normalized Optimization in Deep Neural Networks. *arXiv* (2015), 1–12. arXiv: 1506.02617.
- [71] NEYSHABUR, B., ET AL. Path-Normalized Optimization of Recurrent Neural Networks with ReLU Activations. *arXiv preprint* (May 2016). arXiv: 1605.07154.
- [72] NOVIKOV, A., ET AL. Tensorizing Neural Networks. *Nips* (2015), 1–9. arXiv: arXiv:1509.06569v1.
- [73] OORD, A. van den, ET AL. WaveNet: A Generative Model for Raw Audio (2016), 1–15. arXiv: 1609.03499.
- [74] OORD, A. van den, ET AL. Conditional Image Generation with PixelCNN Decoders. *arXiv* (2016). arXiv: 1606.05328.
- [75] ORS, R. A practical introduction to tensor networks: Matrix product states and projected entangled pair states. *Annals of Physics* 349 (June 2014), 117–158. arXiv: 1306.2164.

- [76] OSEDELETS, I. V. Tensor Train Decomposition. *SIAM J Sci. Comput.* 33, 5 (2011), 2295–2317.
- [77] PASCANU, R., MIKOLOV, T., AND BENGIO, Y. On the difficulty of training recurrent neural networks. *Proceedings of The 30th International Conference on Machine Learning 2* (2012), 1310–1318. arXiv: arXiv:1211.5063v2.
- [78] PASCANU, R., ET AL. How to Construct Deep Recurrent Neural Networks. *CoRR* abs/1312.6 (2013), 1–10. arXiv: 1312.6026.
- [79] PLATE, T. Holographic Reduced Representations. *IEEE Transactions on Neural Networks* 6, 3 (1995), 623–641.
- [80] POLINER, G. E., AND ELLIS, D. P. W. A discriminative model for polyphonic piano transcription. *EURASIP Journal on Advances in Signal Processing* (2007).
- [81] ROSENBLATT, F. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review* 65, 6 (1958), 386–408.
- [82] RUMELHART, D. E., HINTON, G. E., AND WILLIAMS, R. J. Learning representations by back-propagating errors. *Nature* 323, 6088 (1986), 533–536. arXiv: arXiv:1011.1669v3.
- [83] SALIMANS, T. Weight Normalization : A Simple Reparameterization to Accelerate Training of Deep Neural Networks. *arXiv* (2016). arXiv: arXiv:1602.07868v1.
- [84] SAXE, A. M., MCCLELLAND, J. L., AND GANGULI, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *Advances in Neural Information Processing Systems* (Dec. 2013), 1–9. arXiv: 1312.6120.
- [85] SIGAUD, O., ET AL. Gated networks: an inventory. *arXiv* (2015). arXiv: 1512.03201.
- [86] SINGHAL, A. Modern Information Retrieval: A Brief Overview. *Bulletin of the Ieee Computer Society Technical Committee on Data Engineering* 24, 4 (2001), 1–9.
- [87] SMOLENSKY, P. Information processing in dynamical systems: Foundations of harmony theory. *Parallel Distributed Processing Explorations in the Microstructure of Cognition* 1, 1 (1986), 194–281.
- [88] SRIVASTAVA, N., ET AL. Dropout : A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research (JMLR)* 15 (2014), 1929–1958. arXiv: 1102.4807.
- [89] SRIVASTAVA, R. K., GREFF, K., AND SCHMIDHUBER, J. Highway Networks. *arXiv:1505.00387 [cs]* (2015). arXiv: 1505.00387.
- [90] SUTSKEVER, I. Training Recurrent neural Networks. PhD thesis. 2013, 101. arXiv: 1456339 [arXiv:submit].
- [91] SUTSKEVER, I., MARTENS, J., AND HINTON, G. Generating Text with Recurrent Neural Networks. *Neural Networks* 131, 1 (2011), 1017–1024. arXiv: 9809069v1 [arXiv:gr-qc].
- [92] SUTSKEVER, I., ET AL. On the importance of initialization and momentum in deep learning. *Jmlr* 28, 2010 (2013), 1139–1147.
- [93] SZEGEDY, C., IOFFE, S., AND VANHOUCKE, V. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. *Arxiv* (2016), 12. arXiv: 1602.07261.
- [94] TAN, P.-N., STEINBACH, M., AND KUMAR, V. *Introduction to data mining*. Pearson Addison Wesley, 2006, 769.
- [95] TAYLOR, G. W., AND HINTON, G. E. Factored conditional restricted Boltzmann Machines for modeling motion style. In: *Proceedings of the 26th International Conference on Machine Learning (ICML 09)*. 2009, 1025–1032.

- [96] TELGARSKY, M. Benefits of Depth In Neural Networks. In: *29th Annual Conference on Learning Theory*. Vol. 49. 1. 2016, 1–19. arXiv: 1602.04485.
- [97] TENENBAUM, J. B., AND FREEMAN, W. T. Separating style and content with bilinear models. *Neural computation* 12, 6 (2000), 1247–1283.
- [98] THE THEANO DEVELOPMENT TEAM, ET AL. Theano: A Python framework for fast computation of mathematical expressions. *arXiv abs/1605.0* (2016), 19. arXiv: 1605.02688.
- [99] VINCENT, P., ET AL. Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion. *Journal of Machine Learning Research* 11, 3 (2010), 3371–3408. arXiv: 0-387-31073-8.
- [100] VINYALS, O., FORTUNATO, M., AND JAITLEY, N. Pointer Networks. In: *Advances in Neural Information Processing Systems*. 2015. arXiv: arXiv:1506.03134v1.
- [101] VINYALS, O., ET AL. Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 99, PP (2016), 1–1. arXiv: 1609.06647.
- [102] WERBOS, P. J. Backpropagation Through Time: What It Does and How to Do It. *Proceedings of the IEEE* 78, 10 (1990), 1550–1560.
- [103] WU, Y., ET AL. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *arXiv preprint* (2016), 1–23. arXiv: 1609.08144.
- [104] WU, Y., ET AL. On Multiplicative Integration with Recurrent Neural Networks. *arXiv* (June 2016). arXiv: 1606.06630.
- [105] XU, K., ET AL. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. *arXiv preprint* (2015). arXiv: arXiv:1211.5063v2.
- [106] ZAREMBA, W., SUTSKEVER, I., AND VINYALS, O. Recurrent Neural Network Regularization. *arXiv preprint* (2014). arXiv: 1409.2329.
- [107] ZHANG, S., ET AL. Architectural Complexity Measures of Recurrent Neural Networks. In: *ICML*. 2016, 19. arXiv: 1602.08210.