
Understanding RNNs & LSTMs and Reviewing Their Application In Precipitation Nowcasting

Sarwan Shah
Arizona State University
sshah219@asu.edu

Abstract

1 Long Short-Term Memory (LSTM) networks and Recurrent Neural Networks
2 (RNNs) are fundamental architectures in machine learning, particularly when
3 it comes to sequential data analysis. This paper, in an attempt to develop an
4 understanding of these machine learning paradigms, reviews the mathematical
5 theory and modeling of these architectures. It also makes an effort to highlight the
6 practical implications & limitations of these architectures, while also reviewing the
7 literature on their usage in precipitation nowcasting for contextualizing arguments
8 made.

9 1 Introduction

10 The use of modern machine learning, particularly the evolution of deep neural networks, has today
11 opened up a gateway that allows us to capture complex patterns in problems, in data, that have
12 historically been difficult to model mathematically, with an increasingly high degree of accuracy.
13 Consequently, there has been a huge influx of research that is attempting to understand this ability of
14 theirs and also experimenting with their application across a variety of domains. Following this suite,
15 assuming having a foundation established in traditional machine learning modalities (regression,
16 SVM, etc) and their theory, in this paper we attempt to develop an understanding of how modern
17 neural networks work and build upon their traditional counterparts. More particularly our discussion
18 explores the modeling of recurrent neural networks (RNNs) and their variation known as Long
19 Short-Term Memory (LSTM). We highlight some of the key features and limitations of such networks.
20 Following this we also review their application in addressing the challenging problem of precipitation
21 nowcasting. Lastly, we conclude by commenting on the key insights that were taken away from this
22 study.

23 2 Neural Networks

24 2.1 Modelling & Architecture

25 The most basic architecture of a neural network consists of a single neuron that forms a singular
26 hidden layer (in this case) before its output and is fed by inputs via edges that have associated weights,
27 which also form the weights of that neuron. The neuron computes the product of the input and the
28 weights associated with the incoming edge and then applies a function, known as the 'activation
29 function', to give the final output. This architecture is known as a feed-forward network. See figure 1
30 (a).

$$y = \text{sign}\{\mathbf{W}.\mathbf{X}\} = \text{sign}\left\{\sum w_i x_i\right\} \quad (1)$$

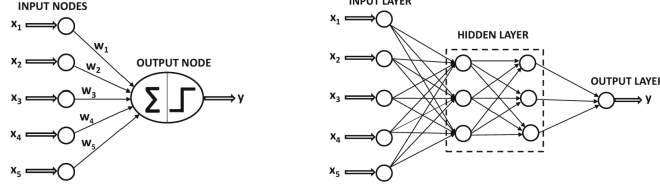


Figure 1: Single (a) & Multi-Layered Neural Network Architecture (b) [1]

It is not difficult to see that with a 'Sign' activation function, this network models the perceptron for binary classification, as shown in equation (1). Similarly, if an identity activation function is used, the network models' linear regression. However, we know such ML algorithms fail to generalize well on data that is not linearly separable.

The additional neurons in each hidden layer and the introduction of more hidden layers - which have inter-connected neurons amongst each other using edges with associated weights - coupled with non-linear activation functions allow for a significant scaling in the expressive power of such networks. Intuitively, with 'Sign' activation it is like combining multiple perceptron-like units to capture statistically complex patterns in data. This is shown in figure 1 (b). And while additional neurons per layer enable capturing more complex patterns, a better trade-off, that has been observed in literature, is to have more layers, or deeper networks, to achieve the same results with fewer overall parameters.

An interesting result that was realized was that the non-linear expressiveness of scaled or deeper neural networks is essentially enabled by the use of non-linear activation functions for neurons, as it was shown in [1] that deep neural networks reduce to linear regression if the neurons consist of only identity activations. Non-linear activations such as 'Sigmoid', 'ReLU', 'Tanh' etc, enable non-linear expressiveness by reducing or mapping the non-linear solution space to a linear one.

Mathematically, it is useful to represent the inputs, outputs, weights, and associated layers in a neural network as vectors and matrices. Thus, for a d_i -dimensional input d_i , and given h_r hidden-layers of length p_r where $r = \{1, 2, \dots, k\}$, the first hidden-layer weights have dimensions $d \times p_1$, where as successive r hidden-layers have weights $p_r \times p_{r+1}$, and the output is $p_k \times d_o$. Additionally, the outputs at each stage of the network are governed by the following equations, where f represents the activation function, \bar{x} the input vector, W the weight matrices, and h_k the output vector of a hidden-layer.

$$\bar{h}_1 = f(W_1^T \bar{x}) \quad (2)$$

$$\bar{h}_r = f(W_r^T \bar{h}_{r-1}) \quad (3)$$

$$\bar{y} = f(W_{out}^T \bar{h}_k) \quad (4)$$

2.2 Training via Back Propagation

Another important aspect of machine learning paradigms is training for a solution that generalizes well over unseen data. The quantity of data required to train neural networks as a rule of thumb is recommended to be 2 to 3 times the number of parameters in the network; which is governed by the number of neurons and edges it has. This ratio helps achieve a good balance against model complexity, helping avoid issues of under or over-fitting. However, this requirement can be relaxed in some cases using techniques such as initialization of weights via pre-training on shallower networks. On the other hand, loss functions used to train neural networks must be differentiable and take the same form as with any optimization-based problem, while their choice remains largely application-specific.

The loss helps optimize edge weights across the network in a stochastic gradient descent fashion using a technique known as backpropagation which calculates the partial derivative of the loss with respect to the weights in the network and updates their weights accordingly. This is enabled by employing the multivariate chain rule in reverse from the output and along the different paths in the network.

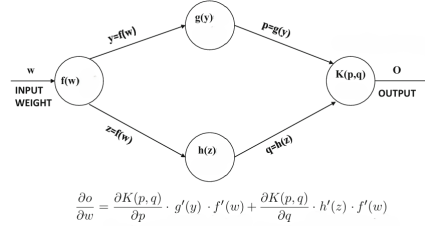


Figure 2: Backpropagation through single-layer two neuron network [1]

2.3 Challenges

A common challenge in the training process for deep neural networks is addressing vanishing and exploding gradients while backpropagating. This happens because the derivatives of activation functions have a saturating nature, in which case, if their gradient computes to < 1 , then, with derivative products taking place amongst multiple layers, the successive gradient quickly approaches zero. On the contrary, if their gradient computes to > 1 , then, their successive gradient can quickly approach larger values.

This inherent instability in training neural networks is partially addressed by choosing activation functions such as 'ReLU' or 'hard Tanh' whose derivatives take values of either 1 (in their region of operation) or 0, and are thus less prone to vanishing or exploding gradients. However, on the flip side, their extremeness can lead to dead neurons (because of 0 gradients if pre-activation values lie outside their region of operation) early on in the training process at times, and thus, additional variants such as 'leaky ReLU' have been devised to address this issue, as they continue to leak some gradient when backpropagating even when pre-activation values lie outside their region of operation. This ensures the neurons don't die entirely suddenly. Additionally, techniques such as batch normalization between layers also help limit the aforementioned along paths in the network.

It is worth commenting that some of this characteristic instability at times is important and adds to the expressiveness of neural networks. For instance, having dead neurons may actually turn out to be essential to be able to learn a desired pattern or relationship in the data.

3 RNNs And LSTMs

3.1 RNNs

3.1.1 Modeling & Architecture

Recurrent Neural Networks are a type of neural network architecture that enables the capturing of temporal relationships in data sequences. In the case of language modeling, the temporal relationship translates to a semantic one between word sequences.

From a modeling perspective, this is achieved by having a separate hidden layer (or layers) in the network for each time-stamp or position, that is input from a time-series or sequence of data. This is done by sharing weights between the consecutive time-stamp hidden layers such that while each time-stamp hidden layer (or layers) has the same input and output weights, it additionally depends on weights from the previous time-stamp hidden layer (or layers), allowing the network to capture temporal dependencies in the data. Architecturally, this is accomplished by incorporating a weighted self-loop, W_{hh} , on the hidden layer which captures the layering through time behavior. This is shown in fig 3 (a) & (b) for a neural network with a single hidden layer. Another key advantage of this architecture is that it allows the network to keep a fixed number of parameters while allowing variable inputs.

Mathematically, in this case, the equations governing the output of the network for $r = 1, 2, \dots, k$ hidden states and a sequence t inputs are (5) and (6), where t captures the layering through time and r through hidden layers. This recursive nature of these equations is what gives RNNs their name.

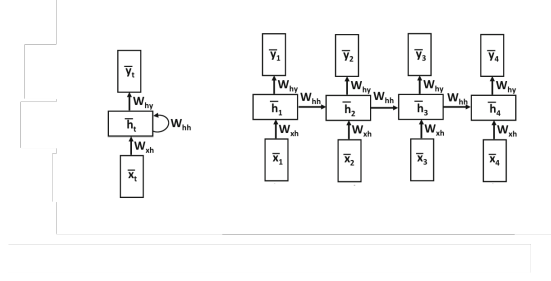


Figure 3: RNN architecture (a) & RNN architecture unfurled through time (b) [1]

$$\bar{h}_t = \tanh(W_{xh}\bar{x}_t + W_{hh}\bar{h}_t - 1) \quad (5)$$

$$\bar{h}_t^r = \tanh(W^r[\bar{h}_t^{r-1}\bar{h}_{t-1}^{r-1}]^T) \quad (6)$$

107 3.1.2 Training via Backpropagation

108 Backpropagation in RNNs similarly takes place as in conventional neural networks. While performing
 109 backpropagation, the layering through time is ignored and the shared weights involved in this process
 110 are assumed to be independent of each other during gradient computations. However, once these
 111 independent gradients have been computed, they are simply summed across the time dimension to
 112 give the actual gradients for W_{xh} , W_{hh} , and W_{hy} . This is known as backpropagation through time
 113 (BPTT).

114 3.1.3 Challenges

115 One of the challenges with training RNNs is that the time dimension can grow significantly long in
 116 some cases. This causes a computational and memory overload during gradient computations. Thus,
 117 a truncated version of BPTT is used, which truncates the gradient computation through time to a
 118 user-defined length.

119 Additionally, owing to the deep depth of such networks during backpropagation updates across the
 120 timed layers, the gradient gets multiplied with the shared weights consistently. If the shared weight w
 121 < 1 , this quickly leads to vanishing gradients, or exploding ones in the opposing case. This adds on
 122 top of the already existing vanishing and exploding gradient problem caused by activation functions.
 123 This makes training RNNs a challenge. To address this inherent instability of RNNs, techniques such
 124 as strong regularization, gradient clipping, and layer normalization have been used. Moreover, it has
 125 given rise to additional architectural variants of the RNNs such as Echo-State networks: randomize
 126 the majority of the hidden-layer weights and only train output layers, LSTMs: maintain a cell state
 127 for stability, and GRUs (Gated Recurrent Units): achieve similar stability as the LSTMs but with
 128 fewer parameters.

129 3.2 LSTMs

130 3.2.1 Modeling & Architecture

131 LSTMs (Long Short-Term Memory) are an architectural variation of RNNs that add more stability by
 132 reducing the vanishing or exploding gradient problem by introducing cell states. These cell states
 133 store information about the past states of the layer and are incorporated as a part of the updates of the
 134 hidden states. Architecturally, this is achieved by introducing an intermediary cell state vector \bar{c}_t^r
 135 that is computed as in (7, 8), via the generation of intermediary input \bar{i} , forget \bar{f} , output \bar{o} , new cell
 136 state \bar{c} gate values. These gate values govern whether the new cell state needs to be added, previously
 137 forgotten/reset, or allowed to leak slightly into the hidden state being updated, as represented by $\bar{i} * \bar{c}$,
 138 $\bar{f} * \bar{c}_t^{r-1}$, and, $\bar{o} * \tanh(\bar{c}_t^r)$ in equations (8, 9) respectively.

$$[\bar{i} \bar{f} \bar{o} \bar{c}]^T = [\text{sigm} \text{sigm} \text{sigm} \text{tanh}]^T(W^r[\bar{h}_t^{r-1}\bar{h}_{t-1}^{r-1}]^T) \quad (7)$$

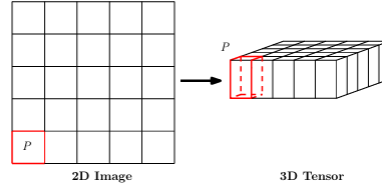


Figure 4: 2D to 3D Tensor [4]

$$\bar{c}_t^r = \bar{f} * c_t^{r-1} + \bar{i} * \bar{c} \quad (8)$$

$$\bar{h}_t^r = \bar{o} * \tanh(\bar{c}_t^r) \quad (9)$$

139 The presence of these cell states adds more randomization and uniqueness to each hidden state
 140 compared to just having the same shared weights temporally. This helps significantly in tackling the
 141 vanishing and exploding gradient problem in conventional RNN architectures.

142 3.3 LSTMs for Precipitation Nowcasting

143 Although a very powerful modeling tool, a natural disadvantage of RNNs and LSTMs is that they
 144 cannot be used with very large sequences of data as in the case of time-series data for weather
 145 forecasting over a period of days. This is because, as explored in the previous sections, very large
 146 sequences lead to a significant increase in the temporal depth of the network, which in turn leads to
 147 computational, memory, and performance-related limitations.

148 However, another realm of weather forecasting known as 'Nowcasting' concerns short-range weather
 149 predictions that usually span a window of 0 to 6 hours. In Meteorology, this is usually considered
 150 a more challenging problem, as it is difficult to model the inherently turbulent nature of weather
 151 phenomena over a short range, and computationally very expensive to simulate. Thus, in this section,
 152 we will explore some literature that experiments with the use of LSTMs to generate short-range
 153 forecasts for the case of precipitation nowcasting.

154 3.3.1 From LSTM to ConvLSTM

155 In [4] they present a neural network architecture that builds upon the LSTM to give a convolutional
 156 LSTM (ConvLSTM) that is used to predict the future states of precipitation radar maps.

157 The key identification they make is that while LSTMs can capture temporal relationships, there
 158 is no way for them to effectively capture spatial relationships, which are often critical to weather
 159 forecasting. To add this spatial awareness they add two spatial dimensions to their inputs, hidden-
 160 states, and outputs, forming them into 3D tensors (see figure 4, thus giving each layer a shape similar
 161 to convolutional layers in a CNN (convolution neural network).

162 Moreover, they use the convolution operation as a part of their recurrence conditions/equation to
 163 ensure spatial patterns are captured, and highlight, how larger kernels can capture the influence of
 164 data points further away, potentially useful for tracking faster movements (e.g., fast-moving weather
 165 systems).

166 They conclude by demonstrating how their ConvLSTM performs significantly better across a variety
 167 of metrics: rainfall mean squared error (Rainfall-MSE), critical success index (CSI), false alarm
 168 rate (FAR), probability of detection (POD), and correlation, compared to state-of-the-art ROVER
 169 algorithm (which is a statistical algorithm for nowcasting using radar maps) and a normal LSTM.

170 3.3.2 From ConvLSTM to TrajGRU

171 This [5] paper builds upon the work conducted in [4]. The key improvements this paper highlights are
 172 that it is used as GRU (Gated Recurrent Unit) instead of the LSTM to introduce performance gains

and reduce the number of parameters required in the network. Moreover, it highlights that even with a ConvGRU, it uses a fixed connection structure and weights for all locations, which may not capture motion patterns that involve rotation and scaling effectively. Thus, it proposes the TrajGRU, which learns a location-variant connection structure that can better model these motion patterns. It does this by calculating continuous optical flows between inputs to represent the connections between locations by embedding it as a part of its recurrence update conditions structure.

The paper concludes by showing how the TrajGRU performs better than ConvLSTM & ConvGRU-based networks using newly defined metrics on the benchmark dataset HKO-07 they created. The new metric adjusts for high-impact cases such as the correct or incorrect prediction of heavy rainfall, and assigning the scenario higher weights.

4 Conclusion

In conclusion, through this paper we have provided a comprehensive exploration into the fundamental architectures of Long Short-Term Memory (LSTM) networks and Recurrent Neural Networks (RNNs), delving into their mathematical theory, modeling, training methodologies, and practical applications.

Key takeaways made through this review were:

an understanding of the evolution of neural networks from single-layer perceptron-like structure to multi-layered architectures capable of capturing complex patterns in data and that is largely enabled by the use of non-linear activation functions. using only linear activations in a neural network reduces them to a linear regression problem. a theoretical and mathematical understanding of why challenges such as vanishing and exploding gradients arise while training deep neural networks, and why the techniques used to mitigate them are effective, such as the use of specific activation functions and normalization techniques. realizing that instabilities in the training process of neural networks add to expressiveness. a theoretical and mathematical understanding of the modeling and architecture of RNNs and LSTMs, how simple modeling tricks enable use to layer sequences through time and preserve states to address critical challenges inherent to the architecture. an appreciation that while RNNs & LSTMs are powerful, their limitations with very long sequences make them less than ideal for tasks like long-term weather forecasting. observing that LSTMs excel at short-term prediction tasks like precipitation nowcasting, where the time window is limited. Understanding the evolution of LSTMs to ConvLSTMs and TrajGRUs to incorporate spatial characteristics to the network and improve their forecasting ability in a targeted manner for the task of precipitation forecasting.

In summary, this paper contributes to the understanding of LSTM and RNN architectures, their training methodologies, and their application in real-world scenarios like weather forecasting. As machine learning continues to advance, further research into enhancing the capabilities of these architectures will undoubtedly lead to even more impactful applications across various domains.

209 **References**

- 210 [1] Aggarwal, C. C. (2023). Neural Networks and deep learning a textbook. Springer International Publishing
211 AG.
- 212 [2] Goodfellow, Ian, et al. “Deep Learning.” Deeplearningbook.org, 2016, www.deeplearningbook.org/. .
- 213 [3] Hochreiter, Sepp, and Jürgen Schmidhuber. “Long Short-Term Memory.” Neural Computation, vol. 9, no. 8,
214 Nov. 1997, pp. 1735–1780.
- 215 [4] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. 2015.
216 Convolutional LSTM Network: a machine learning approach for precipitation nowcasting. In Proceedings of the
217 28th International Conference on Neural Information Processing Systems - Volume 1 (NIPS’15). MIT Press.
- 218 [5] Xingjian Shi, Zhihan Gao, Leonard Lausen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, & Wang-chun Woo.
219 (2017). Deep Learning for Precipitation Nowcasting: A Benchmark and A New Model.
- 220 [6] Hou, L., Zhu, J., Kwok, J., Gao, F., Qin, T., & Liu, T.Y. (2019). Normalization Helps Training of Quantized
221 LSTM. In Advances in Neural Information Processing Systems. Curran Associates, Inc..
- 222 [7] Jiahao Su, Wonmin Byeon, Furong Huang, Jan Kautz, & Animashree Anandkumar (2020). Convolutional
223 Tensor-Train LSTM for Spatio-temporal Learning. CoRR, abs/2002.09131.