



Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

Lecture 4 Sequence to Sequence Learning

Shiwei Lan¹

¹School of Mathematical and Statistical Sciences
Arizona State University

STP598 Advanced Deep Learning Models
Spring 2026



Table of Contents

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

1 Review of RNN

2 Review of LSTM

3 Sequence to Sequence (Seq2Seq) Neural Nets



Recurrent Neural Networks

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

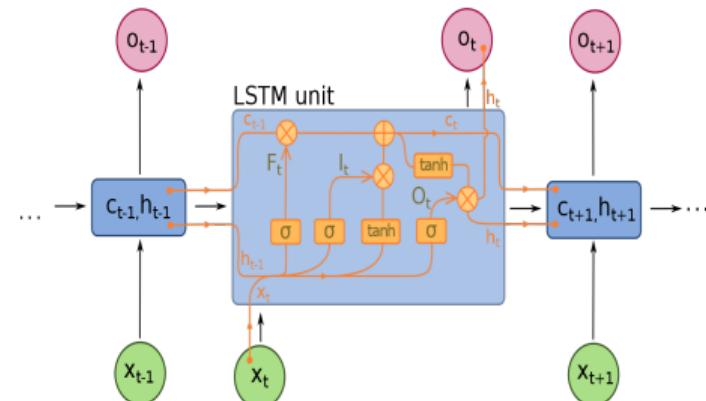
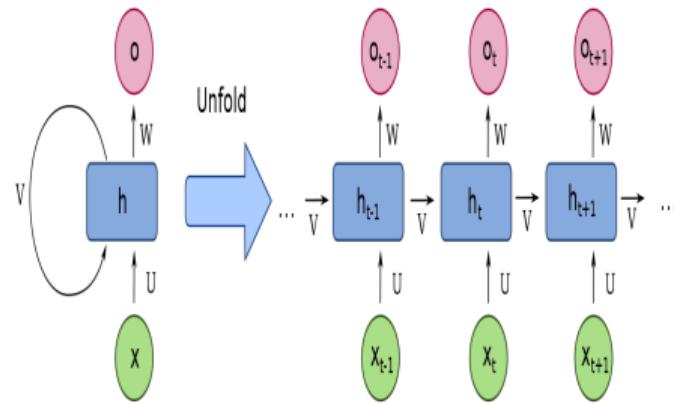
Sequence to
Sequence
(Seq2Seq)
Neural Nets

- Recurrent Neural Networks (RNN) are a family of neural networks for processing sequential data, e.g. time-series, a sequence of images, etc.
- CNN is specialized for processing a grid of values (e.g. images); RNN is specialized for processing a sequence of values (e.g. time series).
- The idea going from MLP to RNN dates back to 1980s: *sharing parameters across different parts of a model.*
- What about CNN across a 1-D temporal sequence? Shallow.
- RNNs designed to retain long time memory, e.g. LSTM, GRU, are particularly useful and successful.

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

Seq2Seq

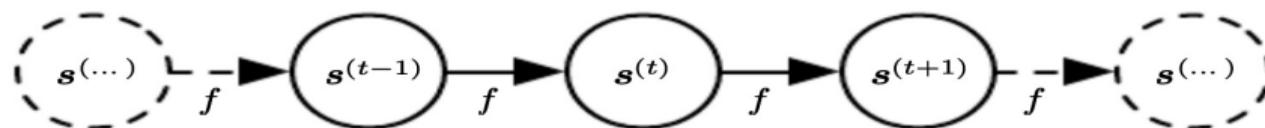
S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

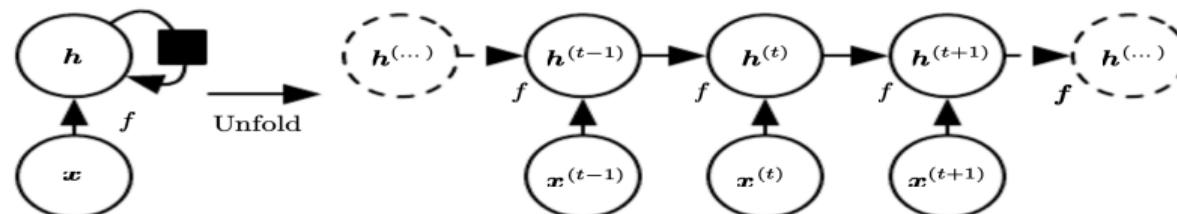
- Consider the classical form of a dynamical system:

$$s^{(t)} = f(s^{(t-1)}; \theta) \quad (1)$$



- Sometimes, the dynamics may be driven by external signal $x^{(t)}$

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta) \quad (2)$$





Unfolding

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- We can represent the unfolded recurrence after t steps with a function $g^{(t)}$:

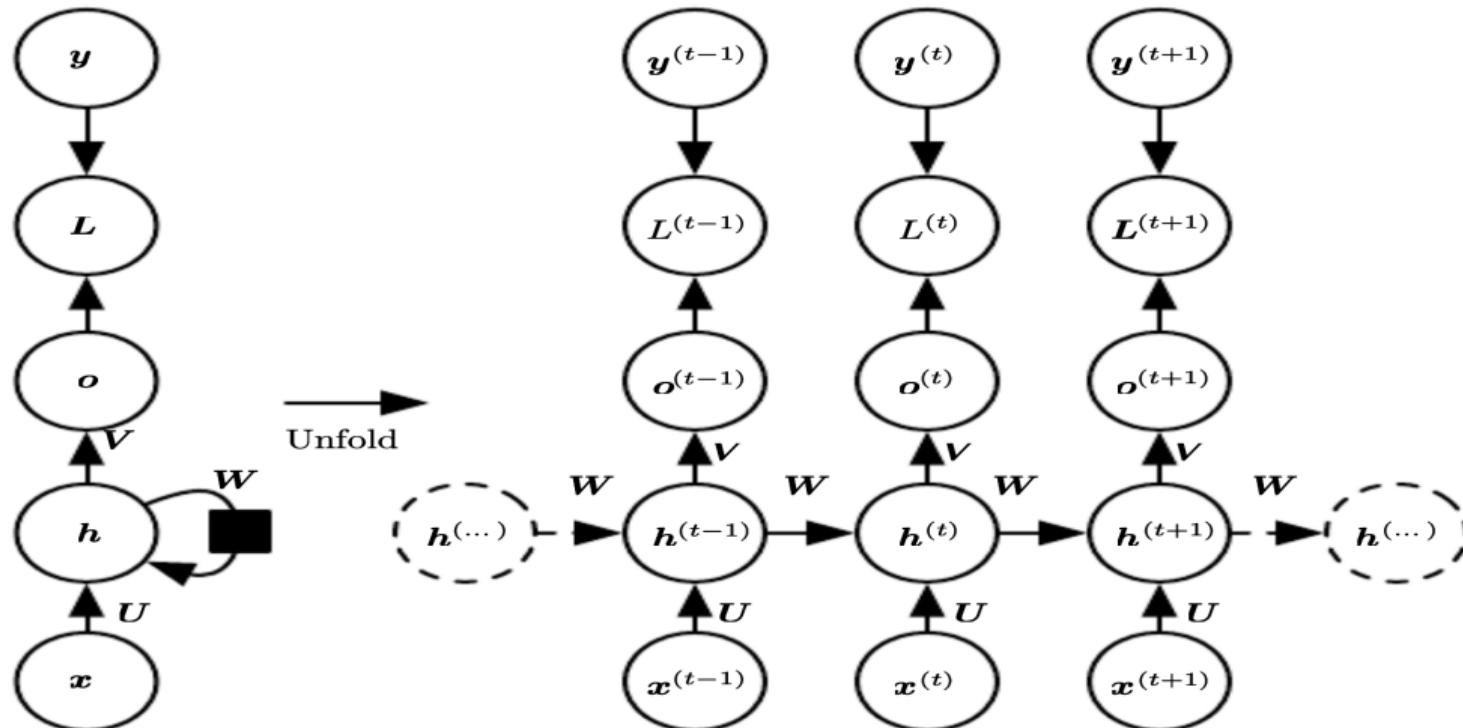
$$h^{(t)} = g^{(t)}(x^{(t)}, x^{(t-1)}, \dots, x^{(1)}) = f(h^{(t-1)}, x^{(t)}; \theta) \quad (3)$$

- The unfolding process thus introduces two major advantages:
 - ① Regardless of the sequence length, the learned model always has the same input size, because it is specified in terms of transition from one state to another state, rather than specified in terms of a variable-length history of states.
 - ② It is possible to use the *same* transition function f with the same parameters at every time step.

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets



Recurrent Neural Networks

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- Assume the output is discrete, as if the RNN is used to predict words or characters.
- Let \hat{y} be a vector of normalized probabilities over the output.
- Forward propagation begins with a specification of the initial state $h^{(0)}$. Then from $t = 1$ to $t = \tau$, we have

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)} \quad (4)$$

$$h^{(t)} = \tanh(a^{(t)}) \quad (5)$$

$$o^{(t)} = c + Vh^{(t)} \quad (6)$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)}) \quad (7)$$

where the parameters are the bias vectors b and c along with the weight matrices U , V and W , respectively for input-to-hidden, hidden-to-output and hidden-to-hidden connections.



Recurrent Neural Networks

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- The total loss for a given sequence of x values paired with a sequence of y values would be the sum of the losses over all the time steps.
- Let $L^{(t)}$ be the negative log-likelihood of $y^{(t)}$ given $x^{(1)}, \dots, x^{(t)}$.

$$\begin{aligned} L(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\}) &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\}) \end{aligned} \tag{8}$$

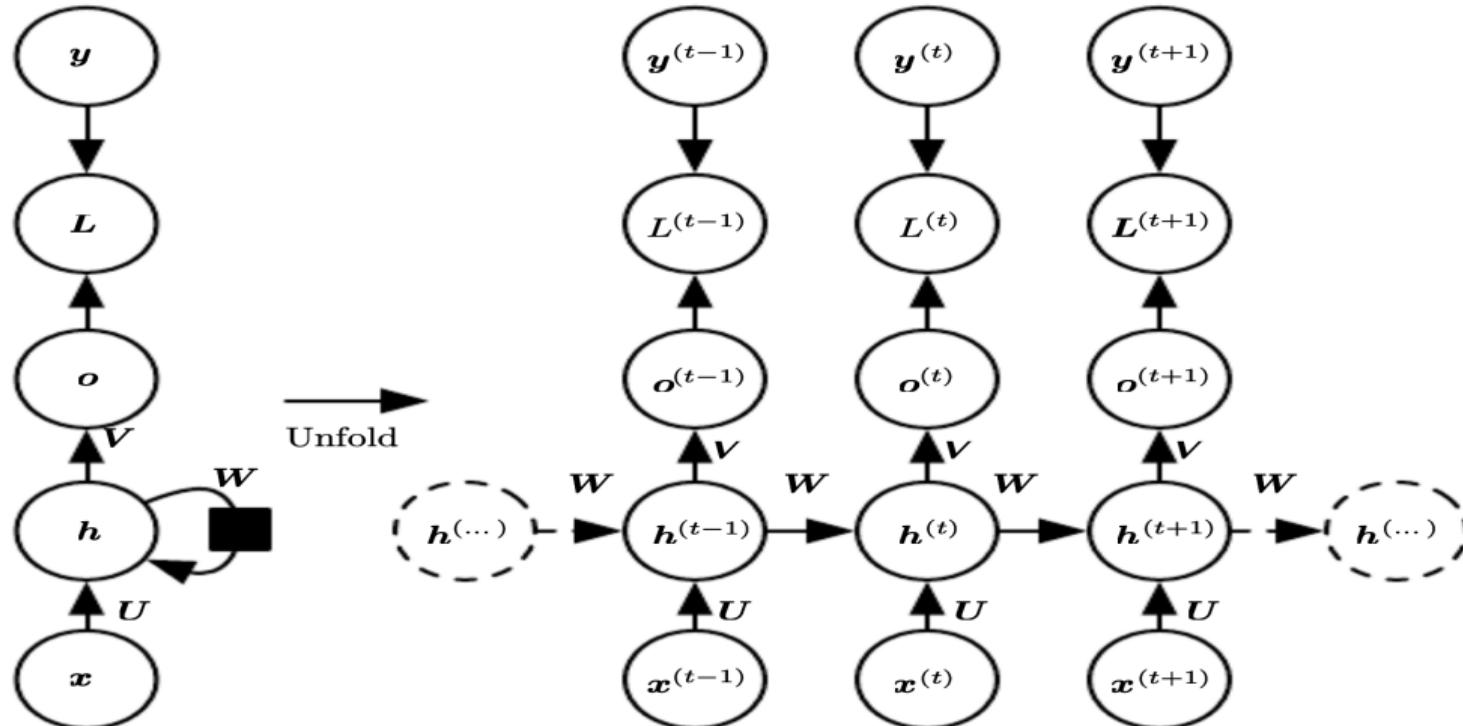
where $p_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\})$ is given by reading the entry for $y^{(t)}$ from the model's output vector $\hat{y}^{(t)}$, i.e. $L^{(t)} = -o^{(t)} + \log \sum_i \exp(o_i^{(t)})$.

- The runtime is $\mathcal{O}(\tau)$; the memory cost is also $\mathcal{O}(\tau)$.
- The back-propagation algorithm applied to the unrolled graph is called **back-propagation through time** (BPTT).

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets



BPTT in RNN

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- We compute the gradients in RNN:

$$\frac{\partial L}{\partial L^{(t)}} = 1 \quad (9)$$

$$(\nabla_{o^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}} \quad (10)$$

$$\nabla_{h^{(\tau)}} L = V^T \nabla_{o^{(\tau)}} L \quad (11)$$

$$\nabla_{h^{(t)}} L = \left(\frac{\partial h^{(t+1)}}{h^{(t)}} \right)^T (\nabla_{h^{(t+1)}} L) + \left(\frac{\partial o^{(t)}}{h^{(t)}} \right)^T (\nabla_{o^{(t)}} L) \quad (12)$$

$$= W^T (\nabla_{h^{(t+1)}} L) \text{diag} \left(1 - (h^{(t+1)})^2 \right) + V^T (\nabla_{o^{(t)}} L) \quad (13)$$



BPTT in RNN

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

$$\nabla_{\mathbf{c}} L = \sum_t \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L \quad (10.22)$$

$$\nabla_{\mathbf{b}} L = \sum_t \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) \nabla_{\mathbf{h}^{(t)}} L \quad (10.23)$$

$$\nabla_{\mathbf{v}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{v}} o_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t) \top} \quad (10.24)$$

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{w}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1) \top} \end{aligned} \quad (10.25) \quad (10.26)$$

$$\begin{aligned} \nabla_{\mathbf{u}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{u}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t) \top} \end{aligned} \quad (10.27) \quad (10.28)$$

Variant: Bidirectional RNNs

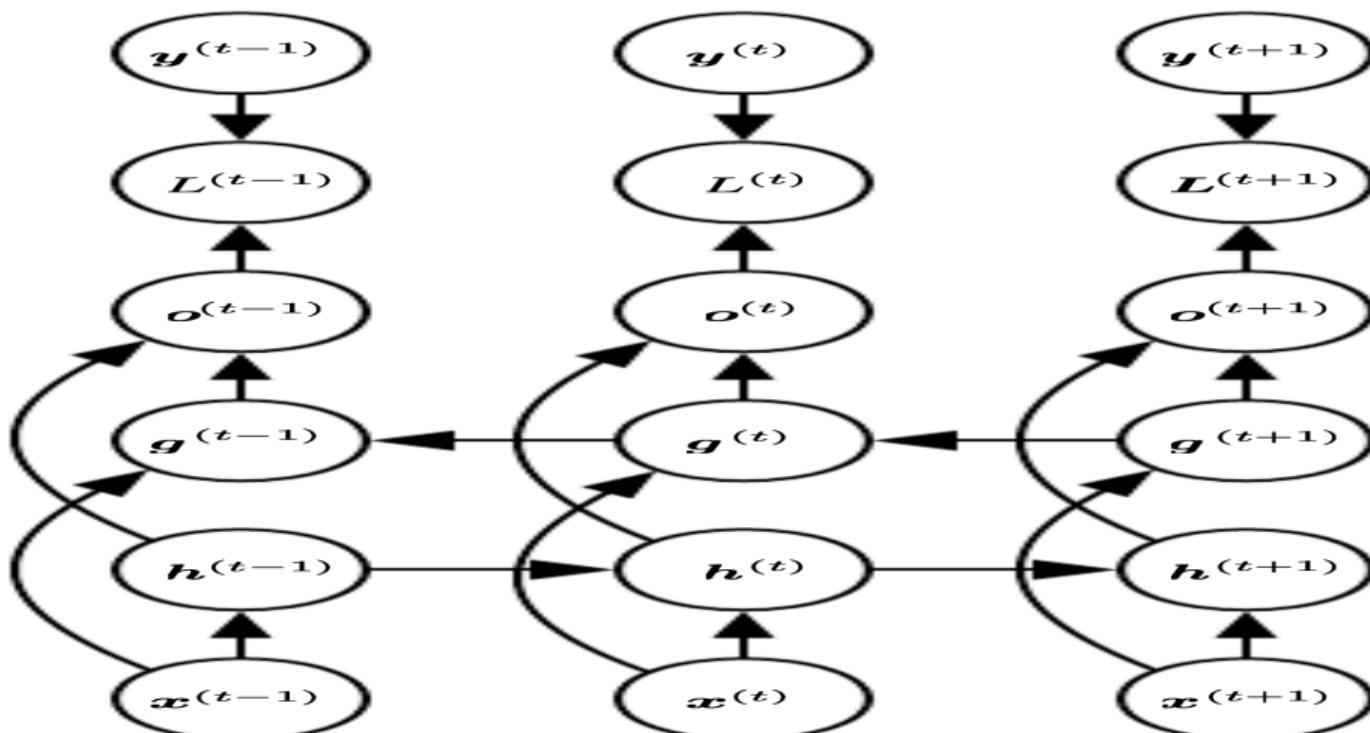
Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets



[Seq2Seq](#)[S.Lan](#)[Review of RNN](#)[Review of LSTM](#)[Sequence to Sequence
\(Seq2Seq\)
Neural Nets](#)

- In many applications we want to output a prediction of $y^{(t)}$ which may depend on *the whole input sequence*.
- Bidirectional RNNs (Schuster and Paliwal, 1997) have been extremely successful (Graves, 2012) in applications such as handwriting recognition (Graves et al., 2008; Graves and Schmidhuber, 2009), speech recognition (Graves and Schmidhuber, 2005; Graves et al., 2013).
- Bidirectional RNNs compute a representation that depends on *both the past (through $g^{(t)}$) and the future (through $h^{(t)}$)* without having to specify a fixed-size window around t .
- This idea can be naturally extended to 2-dimensional input (images) by having four RNNs in 4 directions: up, down, left, right.

What if input and output sequences have different length?

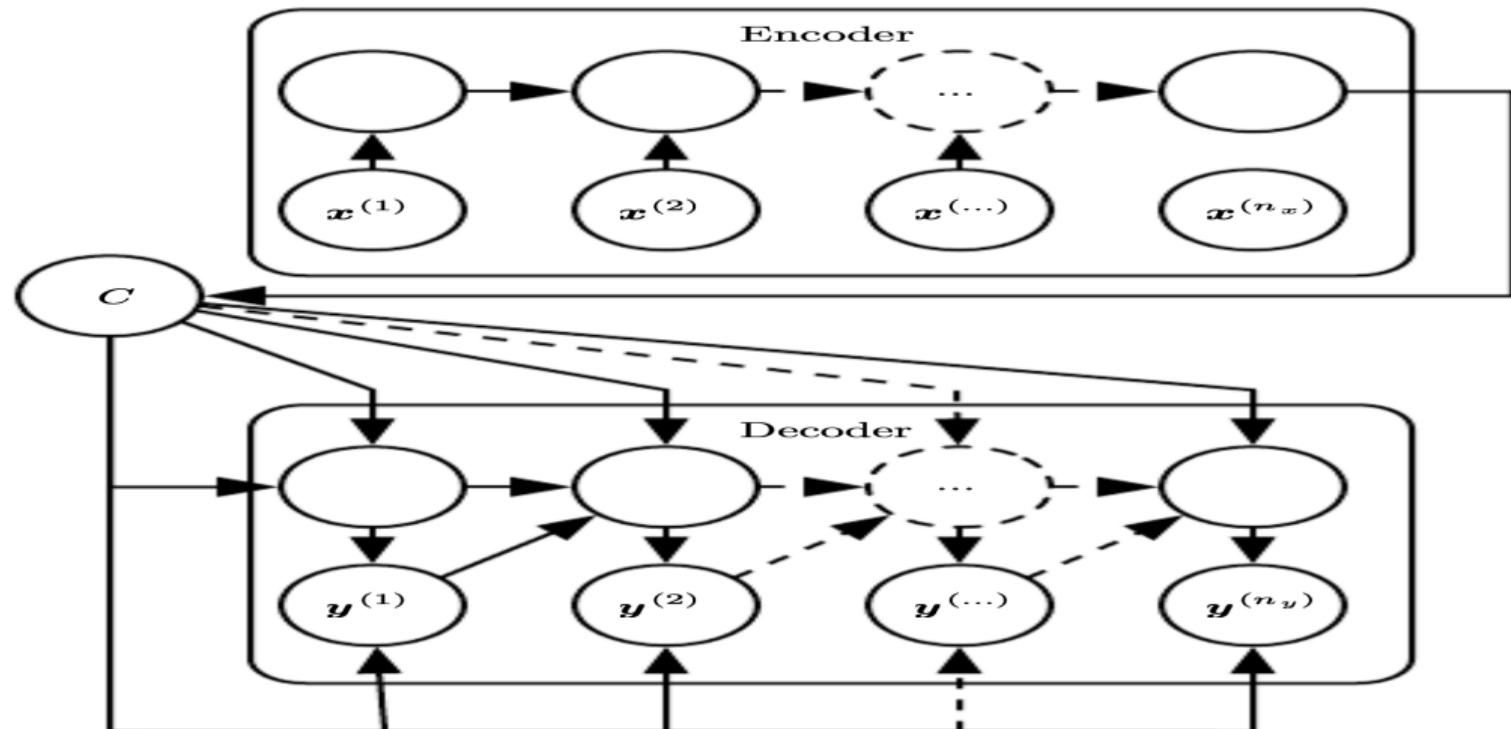
Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets



Variants: Deep RNNs and Recursive NNs

Seq2Seq

S.Lan

Review of RNN

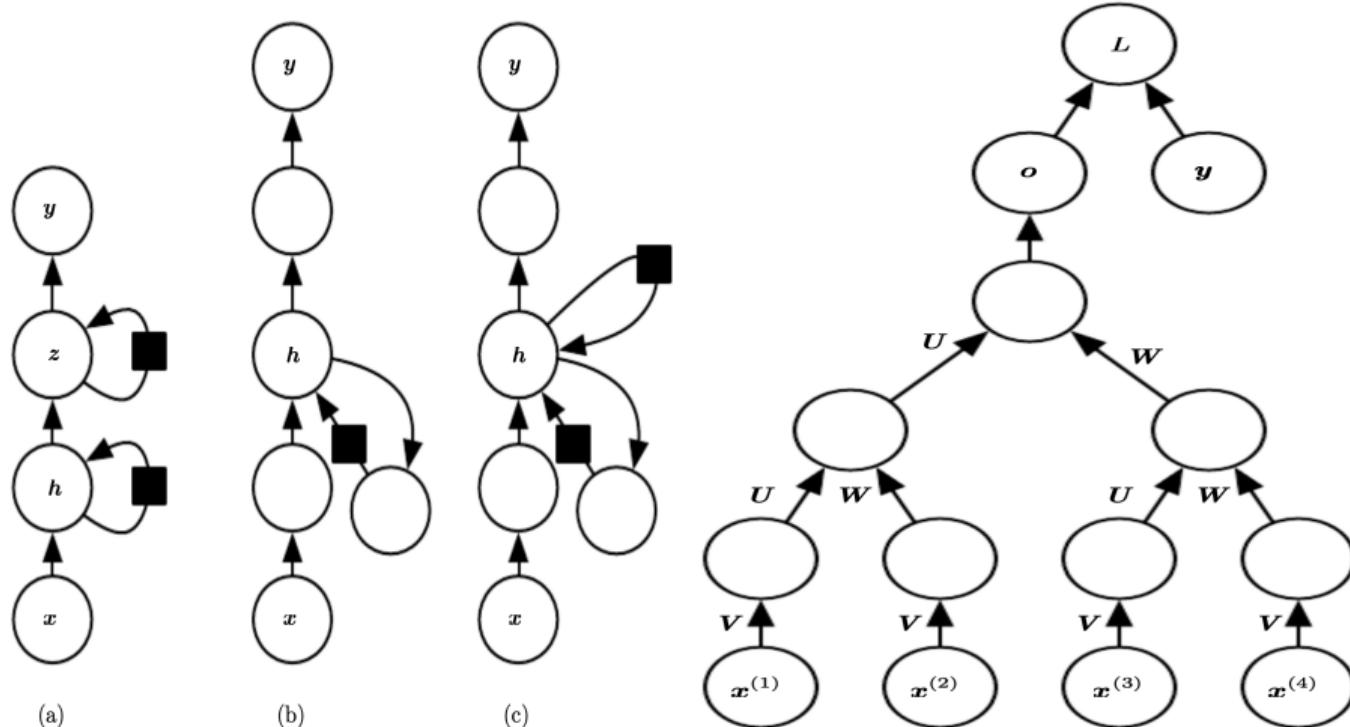
Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets



Table of Contents

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

1 Review of RNN

2 Review of LSTM

3 Sequence to Sequence (Seq2Seq) Neural Nets

The challenge of long-term memory

Seq2Seq

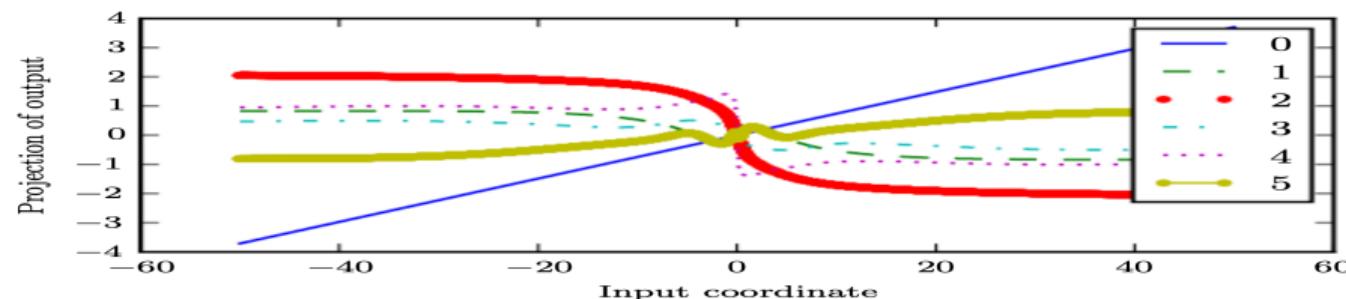
S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)

Neural Nets

- The basic problem is that gradients propagated over many stages tend to either vanish (most of the time) or explode (rarely, but with much damage to the optimization).
- The difficulty with long-term dependencies arises from the exponentially smaller weights given to long-term interactions (involving the multiplication of many Jacobians) compared to short-term ones.



- A heuristic explanation:

$$h^{(t)} = W^T h^{(t-1)} = Q \Lambda Q^T h^{(t-1)} = \dots = Q \Lambda^t Q^T h^{(0)} \quad (14)$$

Long Short-Term Memory (LSTM)

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

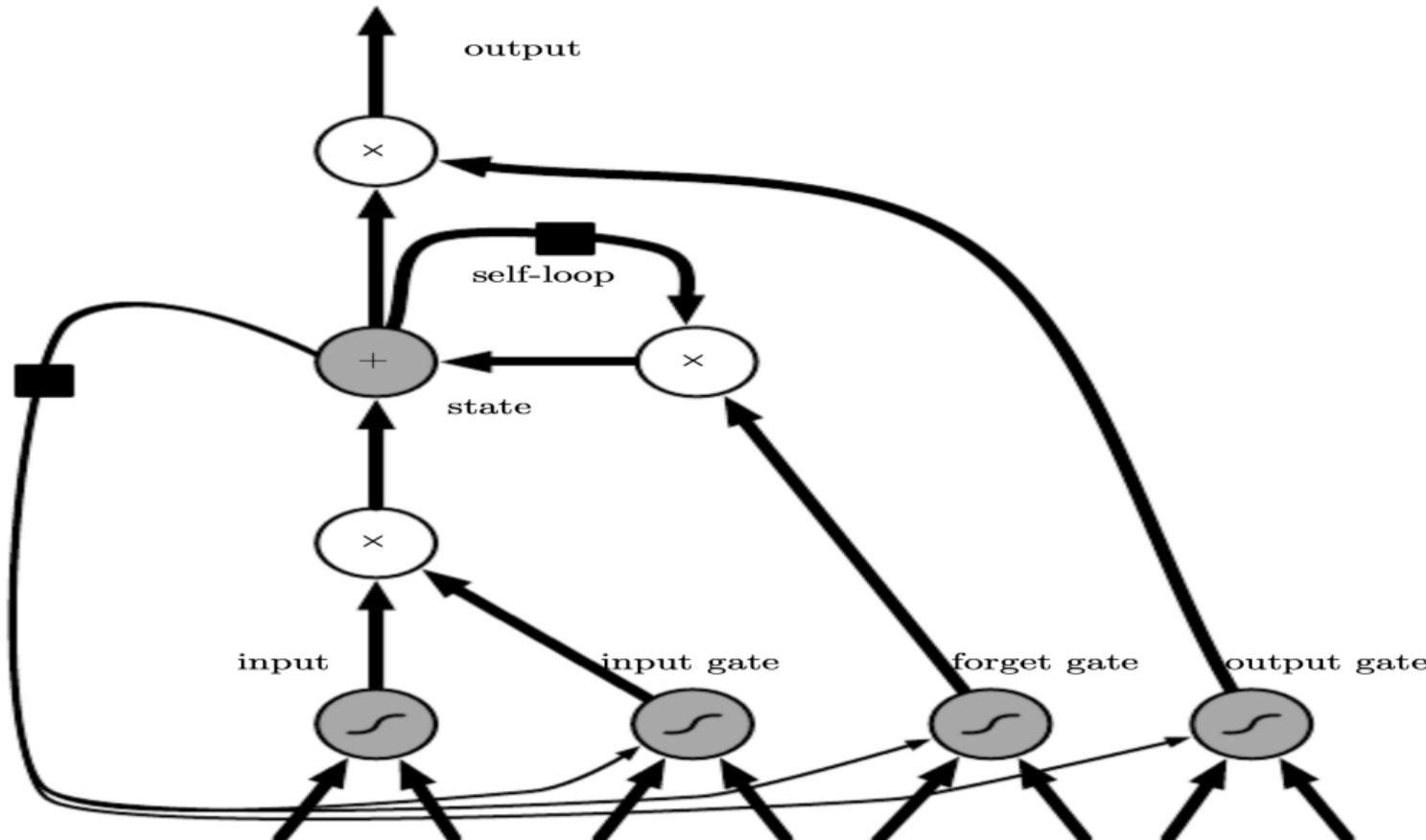
Sequence to
Sequence
(Seq2Seq)
Neural Nets

- The clever idea of introducing self-loops to produce paths where the gradient can flow for long durations is a core contribution of the initial long short-term memory (LSTM) model (Hochreiter and Schmidhuber, 1997)
- A crucial addition has been to make the weight on this self-loop conditioned on the context, rather than fixed (Gers et al., 2000).
- By making the weight of this self-loop gated (controlled by another hidden unit), the time scale of integration can be changed dynamically.
- The LSTM has been found extremely successful in many applications, such as unconstrained handwriting recognition (Graves et al., 2009), speech recognition (Graves et al., 2013; Graves and Jaitly, 2014), handwriting generation (Graves, 2013), machine translation (Sutskever et al., 2014), image captioning (Kiros et al., 2014b; Vinyals et al., 2014b; Xu et al., 2015) and parsing (Vinyals et al., 2014a).

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets



Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

- The self-loop weight (or the associated time constant) is controlled by a **forget gate** unit $f_i^{(t)}$ (for time step t and cell i)

$$f_i^{(t)} = \sigma \left(b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right) \quad (15)$$

- The LSTM cell internal state is updated as follows

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma \left(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right) \quad (16)$$

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

- The **external input gate** unit $g_i^{(t)}$ is computed similarly

$$g_i^{(t)} = \sigma \left(b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right) \quad (17)$$

- The output $h_i^{(t)}$ of the LSTM cell can also be shut off, via the **output gate** $q_i^{(t)}$

$$h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)} \quad (18)$$

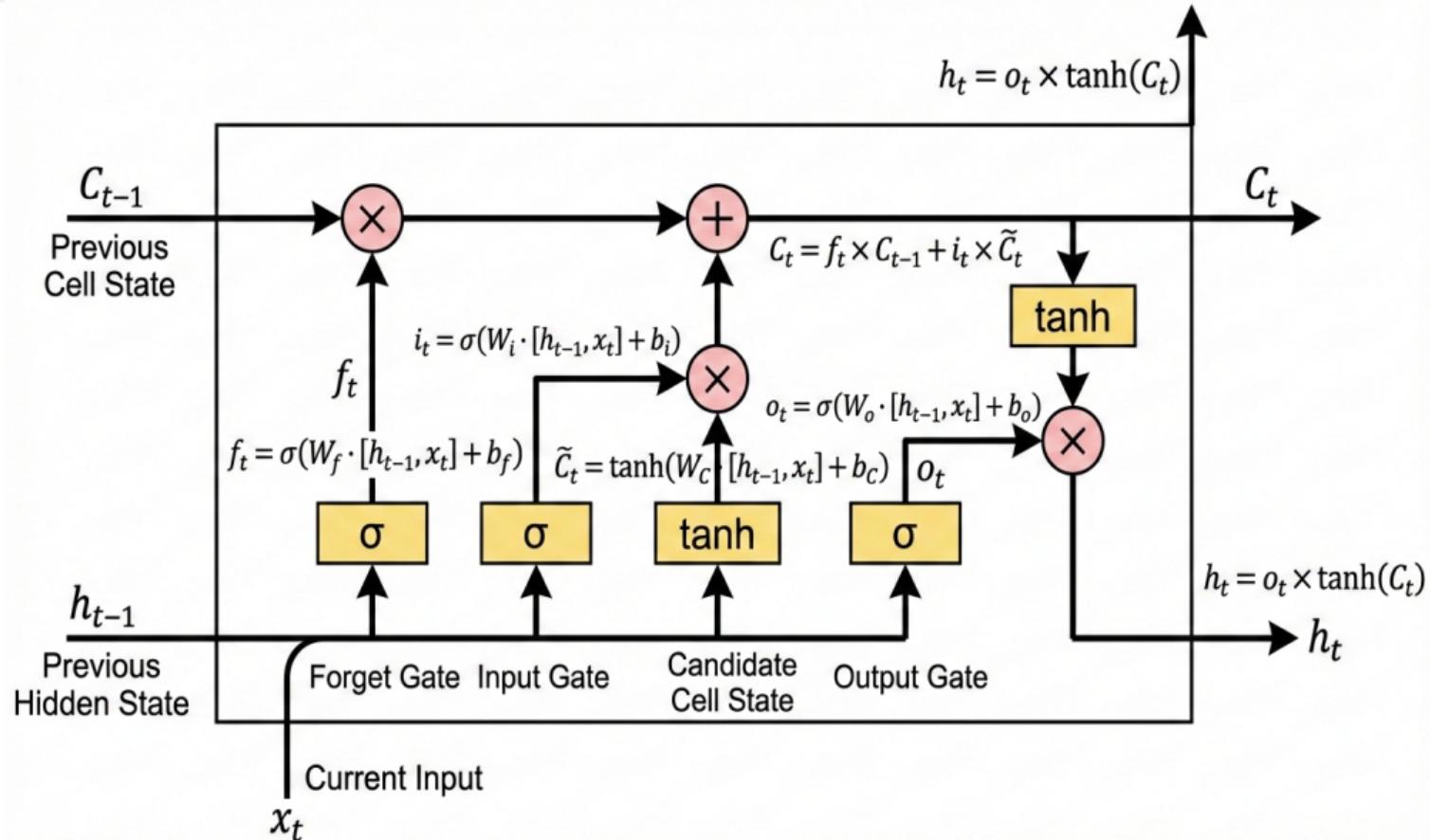
$$q_i^{(t)} = \sigma \left(b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} \right) \quad (19)$$

- A simplified version is called gated recurrent units or GRUs (Cho et al., 2014b; Chung et al., 2014, 2015a; Jozefowicz et al., 2015; Chrupala et al., 2015).

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets



PyTorch Implementation

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- In PyTorch , we build RNN with `torch.nn.RNN`, `torch.nn.LSTM` or `torch.nn.GRU`.
- `torch.nn.RNN` , `torch.nn.LSTM` and `torch.nn.GRU` share some common arguments:
 - `input_size`: the number of expected features in the input x
 - `hidden_size`: the number of features in the hidden state h
 - `num_layers`: number of recurrent layers
 - `nonlinearity` : non-linearity to use, 'tanh' or 'relu'.
 - `bias`: whether to use bias weights $b_i h$ and $b_h h$
 - `bidirectional` : whether to becomes a bidirectional RNN



Table of Contents

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

1 Review of RNN

2 Review of LSTM

3 Sequence to Sequence (Seq2Seq) Neural Nets

[Seq2Seq](#)[S.Lan](#)[Review of RNN](#)[Review of LSTM](#)[Sequence to Sequence \(Seq2Seq\) Neural Nets](#)

- Standard deep neural networks (DNNs) work well in learning labeled data with fixed sizes, not mapping sequences to sequences of variable lengths.
- Sutskever et al (2014) proposed [Seq2Seq](#) by using a multilayered LSTM to map input sequence into a fixed latent space, then another deep LSTM to decode the target sequence from the latent code.
- One important contribution is the *reverse order of input sequence* to feed into encoding LSTM, which shortens the distance of early parts of both sequences, creating short term memory and making training easier.
- Seq2Seq makes minimal assumption on the sequence structure, but can learn long sequences via LSTM.
- It achieved BLEU score of 36.5 in WMT-14 English-French translation, close to the previous SOTA.

Seq2Seq Architecture

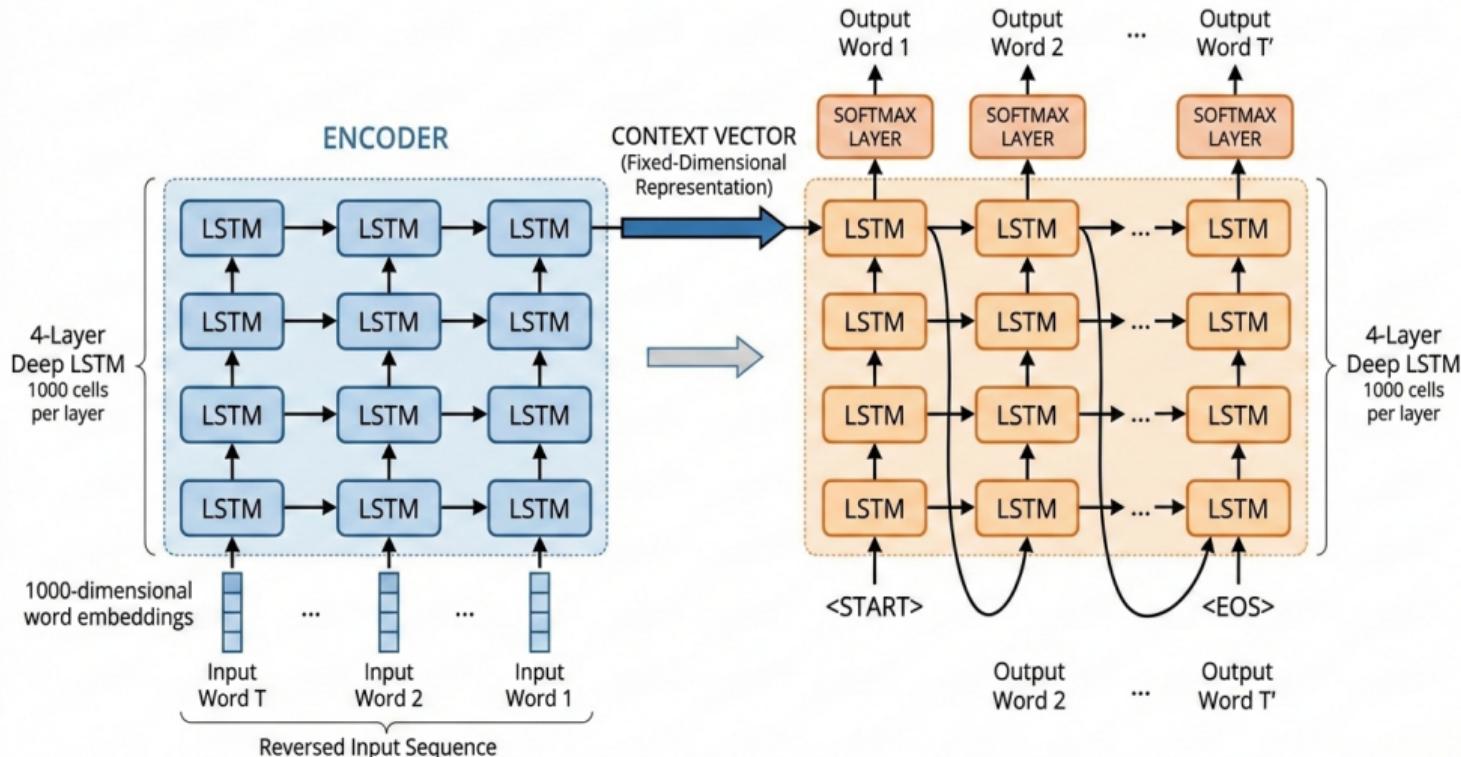
Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

"Sequence to Sequence Learning with Neural Networks"
(Sutskever et al., 2014)





Seq2Seq Model

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- Seq2Seq uses LSTM, a special RNN that maps input sequence (x_1, \dots, x_T) to output sequence (y_1, \dots, y_T)

$$h_t = \sigma(W^{\text{hx}}x_t + W^{\text{hh}}h_{t-1}), \quad y_t = W^{\text{yh}}h_t$$

- LSTM estimates the following conditional probability

$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

where v is the fixed-dimensional latent representation of input sequence (x_1, \dots, x_T) (output of encoder LSTM), and each $p(y_t | v, y_1, \dots, y_{t-1})$ is represented by a softmax over all the words in the vocabulary.

- We train the model by maximizing $1/|\mathcal{S}| \sum_{T,S \in \mathcal{S}} \log p(T|S)$ with input S and output T in the training set \mathcal{S} ; then predict the most probable result by $\hat{T} = \arg \max_T p(T|S)$ over outputs by decoder LSTM.

Reverse the Source

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

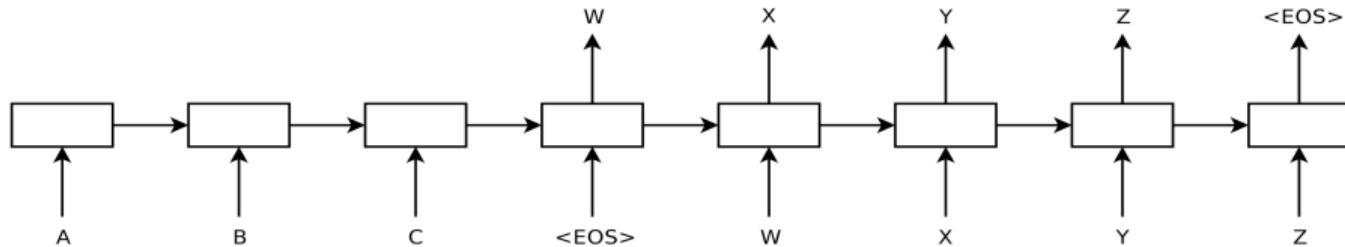


Figure 1: Our model reads an input sentence “ABC” and produces “WXYZ” as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

- Seq2Seq computes the representation of “A”, “B”, “C”, “ EOS_i ” and then uses this representation to compute the probability of “W”, “X”, “Y”, “Z”, “ EOS_i ”.
- With the reverse order of words in the source sentence but not the target sentences, SGD could train easier and LSTM did not suffer on very long sentences.

Numerical Results

Seq2Seq

S.Lan

Review of RNN

Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

Table 1: The performance of the LSTM on WMT'14 English to French test set (ntst14). Note that an ensemble of 5 LSTMs with a beam of size 2 is cheaper than of a single LSTM with a beam of size 12.

Method	test BLEU score (ntst14)
Baseline System [29]	33.30
Cho et al. [5]	34.54
State of the art [9]	37.0
Rescoring the baseline 1000-best with a single forward LSTM	35.61
Rescoring the baseline 1000-best with a single reversed LSTM	35.85
Rescoring the baseline 1000-best with an ensemble of 5 reversed LSTMs	36.5
Oracle Rescoring of the Baseline 1000-best lists	~45

Table 2: Methods that use neural networks together with an SMT system on the WMT'14 English to French test set (ntst14).

Numerical Results

Seq2Seq

S.Lan

Review of RNN

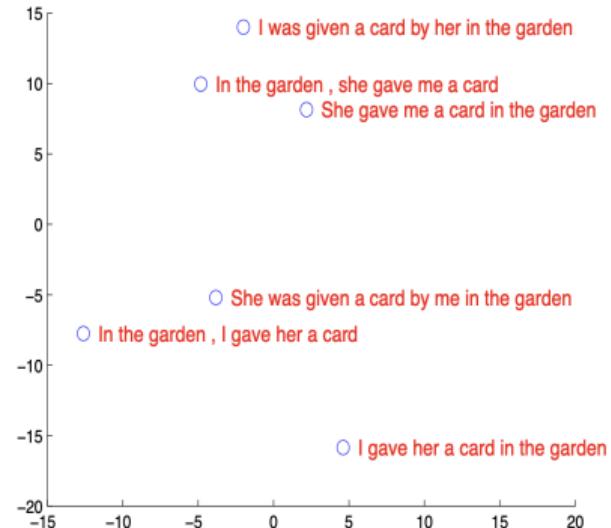
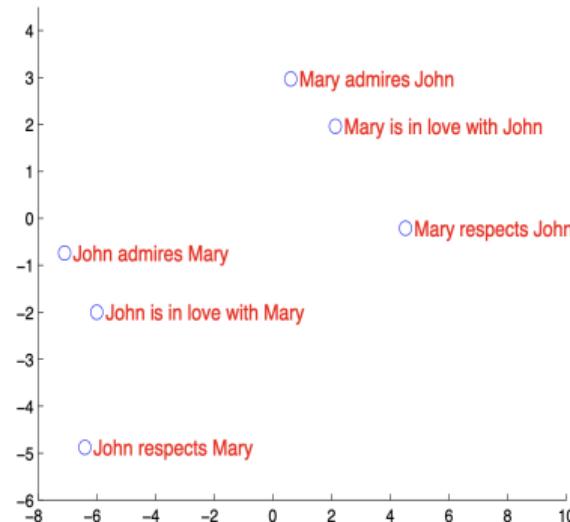
Review of
LSTMSequence to
Sequence
(Seq2Seq)
Neural Nets

Figure 2: The figure shows a 2-dimensional PCA projection of the LSTM hidden states that are obtained after processing the phrases in the figures. The phrases are clustered by meaning, which in these examples is primarily a function of word order, which would be difficult to capture with a bag-of-words model. Notice that both clusters have similar internal structure.



Implementation of Seq2Seq

Seq2Seq

S.Lan

Review of RNN

Review of
LSTM

Sequence to
Sequence
(Seq2Seq)
Neural Nets

- In PyTorch , there is a good tutorial on building Seq2Seq to translate from French to English.
- Here is another Seq2Seq for time series forecasting
<https://discuss.pytorch.org/t/seq2seq-model-with-attention-for-time-series-forecasting/80463>.
- This is a good Github repo for a PyTorch implementation of Seq2Seq
<https://github.com/bentrevett/pytorch-seq2seq>.
- Yet another good tutorial on Seq2Seq for machine translation
<https://machinelearningmastery.com/building-a-plain-seq2seq-model-for-language-translation/>.
- We will demo a Seq2Seq on sol.