

Sequence Modeling: Recurrent Neural Networks

Yagmur Gizem Cinar, Eric Gaussier

AMA, LIG, Univ. Grenoble Alpes

27 October 2017

Deep Learning
Ian Goodfellow and Yoshua Bengio and Aaron Courville
MIT Press
2016

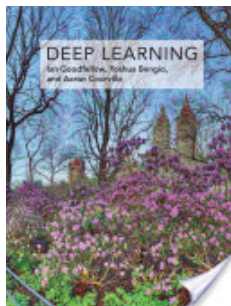


Table of Contents

- 1 Recurrent Neural Networks
- 2 Bidirectional Recurrent Neural Networks
- 3 Encoder-Decoder Sequence-to-Sequence Architectures
- 4 Deep Recurrent Networks
- 5 Long-Term Dependencies
- 6 Leaky Units and Multiple Time Scales
- 7 Long Short-Term Memory and Other Gated RNNs
- 8 Optimization for Long-Term Dependencies

Recurrent Neural Networks (RNNs)

- Specialized for processing a sequence $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}$
- RNNs can scale to longer sequence than networks without sequence-based specialization
- Parameter sharing across different parts of a model enables to extend to different forms and generalize
- For RNNs sharing parameters across time and generalize to different length of sequences
- Example:
 - "I went to Nepal in 2009"
 - "In 2009, I went to Nepal"
- A feedforward network for a fixed sized sentence can learn rules separately at each position
- A RNN share same weights across several time steps

Recurrent Neural Networks (RNNs)

- Parameter sharing with the convolution across 1-D temporal sequence
 - Basis for time-delay networks
 - Parameter sharing across time but shallow
 - Output is a function of neighbouring members
 - Using same convolution kernel at each time step
- For RNNs output is a function of previous members of output
- Output members are produced using the same update rule
- Recurrent parameter sharing leads to a deep computational graph
- A sequence $\mathbf{x}^{(t)}$
 - Time index $t \in [1, \dots, \tau]$
 - In practice, minibatches of sequences with different length τ
 - Time might refer to a position in the sequence

Unfolding Computational Graphs

A computational graph formalizes the structure of a set of computations including mapping inputs and parameters to outputs and loss

unfolding a recurrent/recursive computation to computational graph with a repetitive structure

Classical form of dynamical system:

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

where $\mathbf{s}^{(t)}$ is the state of the system

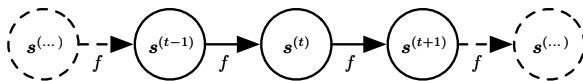


Figure 1: Unfolded computational graph of classical dynamical system¹.

¹ | Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Unfolding Computational Graphs

Classical form of dynamical system:

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

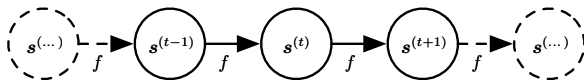


Figure 2: Unfolded computational graph of classical dynamical system².

For finite τ time steps, we can unfold by the same definition $\tau - 1$ times

For $\tau = 3$ time steps:

$$\begin{aligned} \mathbf{s}^{(3)} &= f(\mathbf{s}^{(2)}; \theta) \\ &= f(f(\mathbf{s}^{(1)}; \theta); \theta) \end{aligned}$$

² Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Unfolding Computational Graphs

Another dynamical system driven by an external signal $\mathbf{x}^{(t)}$

$$\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}, \mathbf{x}^{(t)}; \theta)$$

Any function with recurrence can be considered a recurrent network
 Rewriting above equation using variable \mathbf{h} , hidden units

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

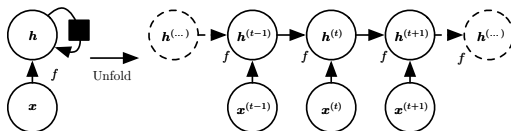


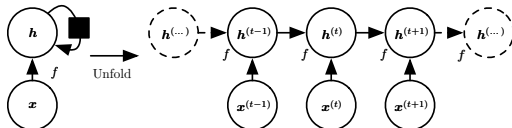
Figure 3: An RNN with no output³.

where has information about the whole sequence

³ | Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

<http://www.deeplearningbook.org>. MIT Press, 2016.

Unfolding Computational Graphs



where

- has information about the whole sequence
- Circuit diagram (left)
- e.g. biological neural network
- Black square indicates a delay of a single time step, from state t to $t + 1$
- Unfolded graph (right) maps a circuit to a computational graph with repeated parts
- Unfolded graph size depends on the sequence length

Unfolding Computational Graphs

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

For a task requiring predicting the future from the past,

- $\mathbf{h}^{(t)}$ becomes a kind of lossy summary
- $\mathbf{h}^{(t)}$ is a fixed-length vector mapping from arbitrary length sequence $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)})$
- Depending on the training criterion, selectively keep some aspects
- Ex: statistical language modeling predict next word given previous words
- Most challenging recovering input sequence from $\mathbf{h}^{(t)}$, e.g. autoencoders

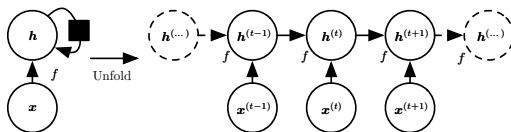
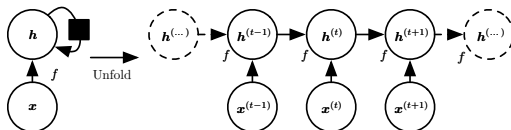


Figure 4: An RNN with no output⁴.

⁴ Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Unfolding Computational Graphs

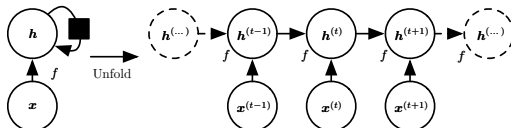


- Unfolded recurrence after t steps with a function $g^{(t)}$

$$\begin{aligned} \mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta) \end{aligned}$$

- $g^{(t)}$ takes whole past sequence $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$ and produce current state
- Unfolded recurrent structure factorize $g^{(t)}$ into repeated f

Unfolding Computational Graphs

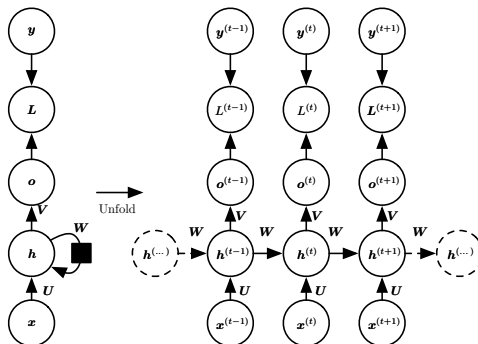


- Unfolding process:
 - + Regardless of sequence length model has fixed input size (since specified in terms of state transitions)
 - + Use of same transition function f with the same parameters at every time step
- These two factors enables to learn a single shared model
 - Generalization to sequence lengths not observed in training
 - Able to train with fewer examples required without than parameter sharing
- Recurrent graph is succinct
- Unfolded graph is explicit and illustrates information flow in forward and backward in time

Recurrent Neural Networks

Important design patterns:

- 1 An output at each time step and recurrent connections between hidden units
- 2 An output at each time step and recurrent connections only from the output at one time step to hidden units at the next time step
- 3 Recurrent connections between hidden units and a single output after reading the entire sequence



Recurrent Neural Networks

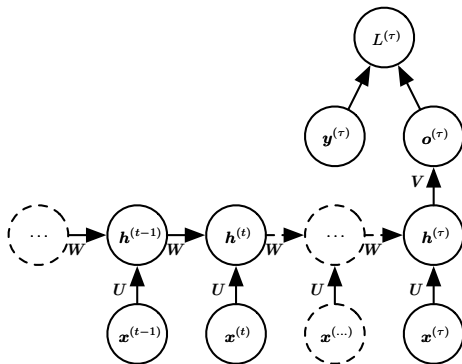


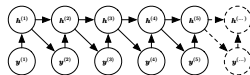
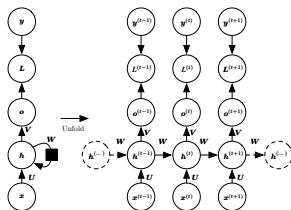
Figure 6: Time-unfolded RNN with a single output at the end⁶.

⁶ | Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

<http://www.deeplearningbook.org>. MIT Press, 2016.

Recurrent Neural Networks

- RNNs below universal for any function computable by a Turing machine
- The output after number of time steps asymptotically linear
 - in the number of time steps used by the Turing machine
 - in the length of the input
- Function computable by Turing machine
 - discrete
 - exact implementation of function (not approximation)
- When an RNN used as a Turing machine
 - input, a binary sequence
 - output, a binary output (discretized)



Recurrent Neural Networks

An RNN outputs at each time step with recurrent connections between hidden units

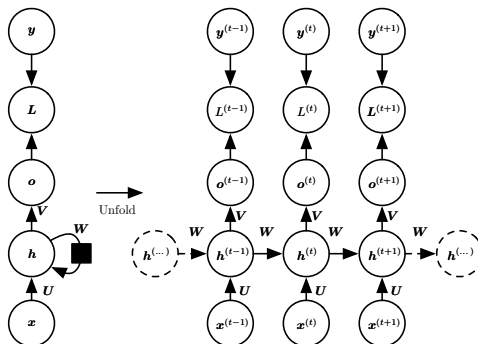


Figure 7: Recurrent hidden units⁷.

⁷I Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Recurrent Neural Networks

An RNN outputs at each time step with recurrent connections between hidden units

$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}$$

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

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}$$

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

\mathbf{b} , \mathbf{c} are biases, weight matrices \mathbf{U} (input-to-hidden), \mathbf{W} (hidden-to-hidden), \mathbf{V} (hidden-to-output)

Total loss between sequence of \mathbf{x} and corresponding sequence of \mathbf{y} :

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

Recurrent Neural Networks

An RNN outputs at each time step with recurrent connections between hidden units

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

$L^{(t)}$ negative log-likelihood of $y^{(t)}$ given $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$

Computing gradient of this loss function expensive

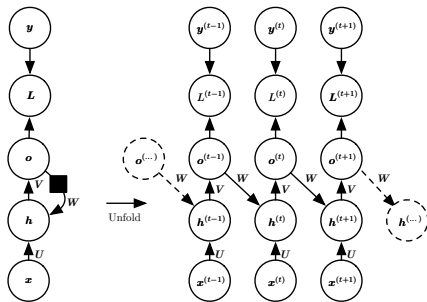
- Forward propagation pass, Backward propagation pass
- Runtime $\mathcal{O}(\tau)$ and cannot be reduced with parallelization
- Memory $\mathcal{O}(\tau)$
- Back-propagation applied to unrolled graph with $\mathcal{O}(\tau)$
back-propagation through time (BPTT)

Recurrent Neural Networks

An RNN outputs at each time step with recurrent connections from output to next step hidden units

- strictly less powerful without hidden-to-hidden recurrence
- cannot simulate a universal Turing machine
- output has to capture all past history
- + training can be parallelized with gradient computed in isolation

For training **Teacher forcing** can be used



Teacher Forcing

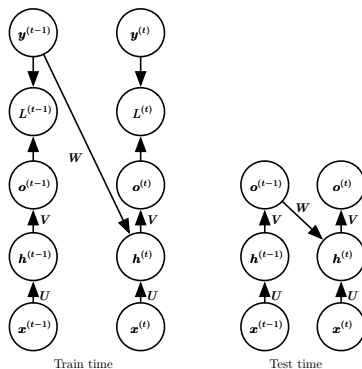


Figure 9: Illustration of teacher forcing⁹.

⁹ | Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

<http://www.deeplearningbook.org>. MIT Press, 2016.

Teacher Forcing

Teacher Forcing emerge from the maximum likelihood criterion

- The conditional maximum likelihood criterion

$$\begin{aligned} \log p\left(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}\right) \\ = \log p\left(\mathbf{y}^{(2)} \mid \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}\right) + \log p\left(\mathbf{y}^{(1)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}\right) \end{aligned}$$

- Enables to avoid BPTT (without hidden-to-hidden connections)
- Teacher forcing can be applied also to models with hidden-to-hidden connections
- But, BPTT is necessary for hidden-to-hidden connections
- Both BPTT and teacher forcing can also be used

Teacher Forcing

Possible disadvantage of strict teacher forcing in use of **open-loop** mode

Open-loop network output feed back as input

Issues with difference in inputs of training and test

- Using both teacher-forced inputs and free-running input during training
- By predicting the correct target a number of steps in the future through unfolded output-to-input paths
- Randomly choose between generated values or actual values during training
- Exploit **curriculum learning** by gradually using more generated values as input

Computing Gradient in RNNs

Computing gradient is applying generalized back propagation to unrolled computation graph

- The computational graph nodes has parameters \mathbf{U} , \mathbf{V} , \mathbf{W} , \mathbf{b} , \mathbf{c} and sequence of nodes indexed by t for $\mathbf{x}^{(t)}$, $\mathbf{h}^{(t)}$, $\mathbf{o}^{(t)}$ and $\mathbf{L}^{(t)}$
- Start with the nodes proceeding the final loss:

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

- Outputs $\mathbf{o}^{(t)}$ input to softmax to obtain $\hat{\mathbf{y}}$ probabilities over output.
- Loss negative log-likelihood of true target $y^{(t)}$ given input
- Gradient $\nabla_{\mathbf{o}^{(t)}} L$ on output at time step t , for all i , t :

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

Computing Gradient in RNNs

At final step τ , $\mathbf{h}^{(\tau)}$ only has \mathbf{o}^τ as a descendent, so gradient:

$$\nabla_{\mathbf{h}^{(\tau)}} L = \mathbf{V}^\top \nabla_{\mathbf{o}^{(\tau)}} L$$

We continue iterating back in time to back-propagate gradients through time, from $t = \tau - 1$ to $t = 1$:

Here, $\mathbf{h}^{(t)}$ (for $t < \tau$) has descendents both \mathbf{o}^t and \mathbf{h}^{t+1}

$$\begin{aligned} \nabla_{\mathbf{h}^{(t)}} L &= \left(\frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{h}^{(t+1)}} L) + \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{o}^{(t)}} L) \\ &= \mathbf{W}^\top (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag} \left(1 - \left(\mathbf{h}^{t+1} \right)^2 \right) + \mathbf{V}^\top (\nabla_{\mathbf{o}^{(t)}} L) \end{aligned}$$

where $\text{diag}(1 - (\mathbf{h}^{t+1})^2)$ represent the diagonal matrix of $1 - (h_i^{t+1})^2$ (Jacobian of hyperbolic tangent of the hidden unit i at time $t + 1$)

Computing Gradient in RNNs

We can obtain the gradients on the parameter nodes once we obtain for internal nodes of computational graph

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

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

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

Computing Gradient in RNNs

We can obtain the gradients on the parameter nodes once we obtain for internal nodes of computational graph

$\nabla_{\mathbf{W}} f$ the contribution of \mathbf{W} to the value of f due to all edges in the graph
 $\mathbf{W}^{(t)}$ are dummy variables as copies of \mathbf{W} but each is used only at time step t
 $\nabla_{\mathbf{W}^{(t)}}$ indicates the contribution of weights at time step t to the gradient

$$\nabla_{\mathbf{W}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_i^{(t)} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top},$$

$$\nabla_{\mathbf{U}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}} h_i^{(t)} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top}$$

No need to compute gradient with respect to $\mathbf{x}^{(t)}$ since no parameter as ancestors

RNNs as Directed Graphical Models

Loss should be chosen according to given task

Interpret output as a probability distribution

Define loss function using the cross-entropy associated with this distribution

- An RNN to estimate the conditional distribution of the next sequence element $\mathbf{y}^{(t)}$ given past elements, we maximize log-likelihood:

$$\log p(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$

- where outputs \mathbf{y} conditionally independent given the sequence \mathbf{x}
- If there exist connections from output at one time step to next time step

$$\log p(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}, \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t-1)})$$

- An RNN models scalar random variables $\mathbb{Y} = \{y^{(1)}, \dots, y^{(\tau)}\}$ with no additional inputs (input at t is output at $t - 1$), the joint distribution for conditional probabilities:

$$P(\mathbb{Y}) = P(y^{(1)}, \dots, y^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} | y^{(t-1)}, y^{(t-2)}, \dots, y^{(1)})$$

RNNs as Directed Graphical Models

- An RNN models scalar random variables $\mathbb{Y} = \{y^{(1)}, \dots, y^{(\tau)}\}$ with no additional inputs (input at t is output at $t - 1$), the joint distribution for conditional probabilities:

$$P(\mathbb{Y}) = P(y^{(1)}, \dots, y^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} | y^{(t-1)}, y^{(t-2)}, \dots, y^{(1)})$$

- Negative log-likelihood:

$$L = \sum_t L^{(t)}$$

where $L^{(t)} = -\log P(y^{(t)} | y^{(t-1)}, y^{(t-2)}, \dots, y^{(1)})$

RNNs as Directed Graphical Models

The edges of the graph represent direct dependency

Markov assumption edges from $\{y^{(t-k)}, \dots, y^{t-1}\}$ to y^t rather than whole history

The graphical model over y values with complete graph structure:

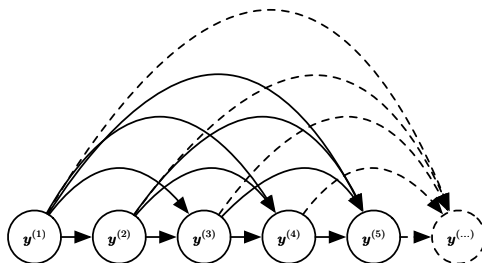


Figure 10: Fully connected graphical model¹⁰.

where $\mathbf{h}^{(t)}$ is ignored

¹⁰ Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

RNNs as Directed Graphical Models

Provides efficient parametrization over observations

If y can take k different values, tabular representation $\mathcal{O}(k^\tau)$

RNN uses the same function f, θ at each time step

- the number of parameters in RNN $\mathcal{O}(1)$ as a function of sequence length
- Due to reduced number of parameters optimizing might be difficult
- Assumption of stationarity, the conditional probability over variables at $t + 1$ given the variables at time t is stationary

With $\mathbf{h}^{(t)}$ nodes decouples past and future

E.g. A variable $y^{(i)}$ may influence $y^{(t)}$ through \mathbf{h}

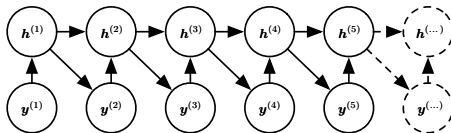


Figure 11: Conditional distributions for hidden units are deterministic¹¹.

¹¹ Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Drawing samples from an RNN model

Sampling from conditional distribution at each time step

How to determine the length of the sequence:

- If output is a symbol from a vocabulary, having a special symbol indicate the end of a sequence
- An extra Bernoulli output to decide to continue or halt
 - E.g. RNN produces a sequence of real numbers, new output unit usually sigmoid with cross entropy loss.
 - Sigmoid maximize log-probability of the sequence ends or not
- An extra output to determine the sequence length τ
 - An extra output predicts the integer τ
 - E.g. sample an integer τ and then sample τ steps of data
 - Here, RNN needs an extra input consist of value of τ or number of remaining steps $\tau - t$
 - Extra input to avoid abrupt ending sequence

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}) = P(\tau)P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)} | \tau)$$

Modeling Sequences Conditioned on Context with RNNs

RNNs with a single vector of \mathbf{x} as input instead of sequence of vectors $\mathbf{x}^{(t)}$
 Providing an extra input to an RNN:

- as an extra input at each time step
- as the initial state $\mathbf{h}^{(0)}$
- both

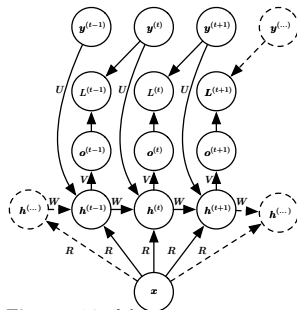


Figure 12: Vector to sequence¹².

¹²| Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

Modeling Sequences Conditioned on Context with RNNs

An input sequence of $\mathbf{x}^{(t)}$ instead of a single input

Conditional distribution of $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)})$ makes a conditional independence assumption

$$\prod_t P(\mathbf{y}^t | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$

By adding connections from output at time t to hidden unit at time $t + 1$:
Hence, the output values are not forced to be conditionally independent for this model:

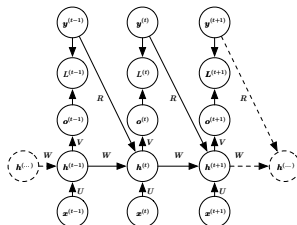


Figure 13: Hidden and output recurrence¹³

Bidirectional RNNs

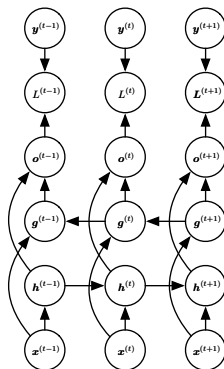


Figure 14: Bidirectional RNN¹⁴.

¹⁴| Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

<http://www.deeplearningbook.org>. MIT Press, 2016.

Bidirectional RNNs

- A causal structure, state at time t captures only past
- Many applications output $\mathbf{y}^{(t)}$ after processing the whole input sequence
e.g. in speech recognition, handwriting recognition
- Bidirectional RNN combines RNN moves forward through time \mathbf{h} and backward through time \mathbf{g}

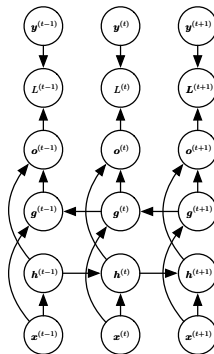


Figure 15: Bidirectional RNN^a.

^aI Goodfellow, Y Bengio, and A Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.

Encoder-Decoder Sequence-to-Sequence Architectures

An RNN map an input sequence to an output sequence which is not necessarily the same length

e.g. speech recognition, machine translation, question answering

The input to RNN is called the context C , summarize the input sequence

$$\mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_x)})$$

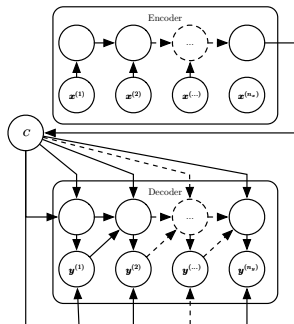


Figure 16: Sequence-to-sequence Architecture¹⁵.

Encoder-Decoder Sequence-to-Sequence Architectures

An **encoder** or **reader** or **input** process input to the sequence

A **decoder** or **writer** or **output** conditioned on a fixed length vector to generate $Y = (y^{(1)}, \dots, y^{(n_y)})$ Two RNNs trained to jointly maximize average of $\log P(y^{(1)}, \dots, y^{(n_y)} | x^{(1)}, \dots, x^{(n_x)})$

The last state of encoder h_{n_x} typically used as C

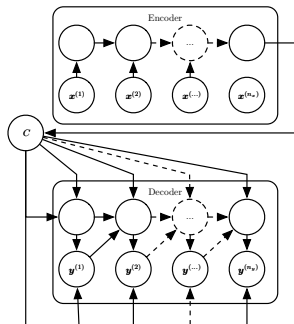


Figure 17: Sequence-to-sequence Architecture¹⁶.

Deep Recurrent Networks

RNNs decomposed into three main blocks of parameters and associated transformations:

- from input to hidden state
- from previous hidden state to next hidden state
- from hidden state to output

Increasing the depth of the RNNs improves

- (a) hierarchical hidden recurrent states
- (b) deeper computation introduced to input-to-hidden, hidden-to-hidden, and hidden-to-output.
- (c) Skip connections can handle path-lengthening effect

Deep Recurrent network larger capacity of representation but might increase the difficulty of optimization

Deep Recurrent Networks

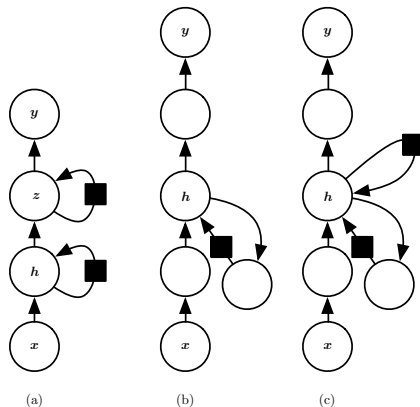


Figure 18: An RNN can be deep many ways. (a) The recurrent state organized into groups hierarchically. (b) Introduced in input-to-hidden, hidden-to-hidden, and hidden-to-output. (c) Skip connections can handle path-lengthening effect¹⁸.

¹⁷ | Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

The Challenge of Long-Term Dependencies

RNNS might construct a very deep computational graphs by repeatedly applying the same operation at each time step of a long sequence

Gradients propagated many stages tends to vanish or explode

Even if the network is stable (can store memories and gradients not exploding), exponentially smaller weights are given to long-term interactions than short ones
 Recurrence relation (for a simple network without nonlinear activation and input \mathbf{x})

$$\mathbf{h}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)}$$

After t time steps (repeatedly multiplying with \mathbf{W})

$$\mathbf{h}^{(t)} = (\mathbf{W}^t)^T \mathbf{h}^{(0)}$$

\mathbf{W} has an eigendecomposition $\mathbf{W} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$, with orthogonal \mathbf{Q}

$$\mathbf{h}^{(t)} = \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)}$$

Leaky Units and Multiple Time Scales

One way of dealing with long term dependencies having multiple time scales
By combining

- Fine-grained time scales
- Coarse time scales

Efficiently transfer information from distant past to present
Various strategies for building both fine and coarse time scales

- Skip connections across time
- Leaky units integrate signals with different time constants
- Removal of some connections

Skip Connections Through Time

Gradients may vanish or explode with respect to number of time steps

- Recurrent connections with time delay of d
- Gradients diminish exponentially as a function of $\frac{\tau}{d}$ rather than τ
- Gradient might still explode with delayed and single step connections
- Skip connections allow to capture longer dependency
- But they might not represent all long-term dependencies

Leaky Units and Spectrum of Different Time Scales

Another way to obtain paths on which the product of derivatives close to 1

- Having units with linear self-connections and a weight near 1 on these connections
- Accumulate a running average $\mu^{(t)}$ of some value $\nu^{(t)}$
- By applying update rule $\mu^{(t)} \leftarrow \alpha \mu^{(t-1)} + (1 - \alpha) \nu^{(t)}$
- α parameter is an example of linear self connection from $\mu^{(t-1)}$ to $\mu^{(t)}$
- When α is near 1, the running average remembers information about past for a long time
- When α is near 0, information about past rapidly discarded
- Hidden units with linear self-connections can behave similar to running averages
- Those hidden units called **leaky units**

Leaky Units and Spectrum of Different Time Scales

Linear self-connection with a weight near 1 is a way of ensuring the access values from past

- More smoothly and flexibly compared to skip connections
- By adjusting α rather than skip-length d

Two strategy of setting time constants of leaky units:

- Fixing values manually remain constant
- As free parameters and learn them

Removing Connections

Another approach to dealing with long-term dependencies is organizing the state of an RNN at multiple scales

- Information flowing easily at a slower time scales
- Actively removing single step connections and replacing them with longer connections
- Units forced to operate on a long time scale
- Compared to skip connections through time
 - Skip connections add edges
 - Skip connections can learn to operate on a long time scale or focus on short-term connections

Gated RNNs

The most effective sequence models used in practical applications are **gated RNNs**

- **Long short-term memory (LSTM)**
- **gated recurrent unit (GRU)**

Like leaky units goal is to create paths through time that have derivatives do not vanish nor explode

- Leaky units has connection weights that manually chosen or learned
- Gated RNNs generalizes this to connections weights that may change at each time step
- Leaky units allows to accumulate information
- But once this information is used, it might be useful to forget
- A mechanism to forget the old state by setting it to 0
- Gated RNNs learn to decide when to forget a state

Long Short-Term Memory (LSTM)

- Core contribution of initial LSTM model¹⁹ is self-loops to introduce paths that gradient can flow
- Gers et al. make the weight on this self-loop weight conditioned on the context²⁰
- With gated weight of self-loop (controlled by another unit), time scale integration dynamically controlled
- The LSTM is very successful in many domains: handwriting detection and generation, time series forecasting, machine translation, speech recognition..
- LSTM recurrent networks have "LSTM cells" that have internal recurrence, a self loop, in addition to outer recurrence of RNN

¹⁹Sepp Hochreiter and Jürgen Schmidhuber. "Long Short-Term Memory". In: *Neural Comput.* 9.8 (Nov. 1997), pp. 1735–1780.

²⁰Felix A. Gers, Jürgen A. Schmidhuber, and Fred A. Cummins. "Learning to Forget: Continual Prediction with LSTM". In: *Neural Comput.* 12.10 (Oct. 2000), pp. 2451–2471. 

Long Short-Term Memory (LSTM)

- Each LSTM cell has the same inputs and outputs with gating units controlling information flow
- State unit $s_i^{(t)}$ has a linear self-loop
- Self-loop weight of $s_i^{(t)}$ controlled by a **forget gate** unit $f_i^{(t)}$ for time step t and cell i
- $f_i^{(t)}$ sets weight of $s_i^{(t)}$ a value of (0,1) via sigmoid unit:

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

- $\mathbf{x}^{(t)}$ is the current input, $\mathbf{h}^{(t)}$ is the current hidden layer vector
- \mathbf{b}^f bias, \mathbf{U}^f input weights and \mathbf{W}^f recurrent weights for the forget gates

Long Short-Term Memory (LSTM)

- The LSTM cell internal state update:

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

- b** bias, **U** input weights and **W** recurrent weights into the LSTM cell
- The **external input gate** $g_i^{(t)}$ is calculated as forget gate but with own parameters:

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

- The output $h_i^{(t)}$ of the LSTM cell can also be shut off via **output gate** $q_i^{(t)}$ which also uses a sigmoid unit for gating:

$$h_i^{(t)} = \tanh \left(s_i^{(t)} \right) q_i^{(t)}$$

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

Long Short-Term Memory (LSTM)

Variant of LSTM uses the cell state $s_i^{(t)}$ as an extra input to the three gates, this would require three additional parameters

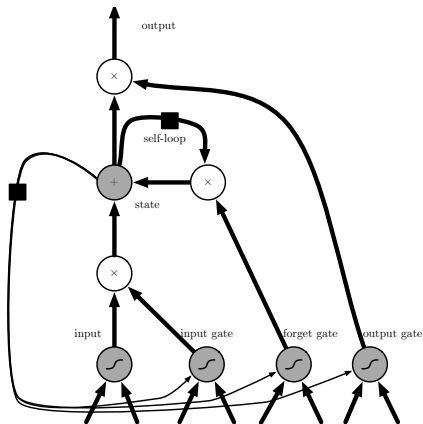


Figure 19: Block diagram of the LSTM²¹.

Gated Recurrent Unit (GRU)

Comparison to LSTM there is a single gating unit simultaneously controls the forgetting factor and the decision to update the state unit

- The update equations

$$h_i^{(t)} = u_i^{(t-1)} h_i^{(t-1)} + (1 - u_i^{(t-1)}) \sigma \left(b_i + \sum_j U_{i,j} x_j^{(t-1)} + \sum_j W_{i,j} r_j^{(t-1)} h_j^{(t-1)} \right)$$

- **u** stands for update gate:

$$u_i^{(t)} = \sigma \left(b_i^u + \sum_j U_{i,j}^u x_j^{(t)} + \sum_j W_{i,j}^u h_j^{(t)} \right)$$

- **r** stands for reset gate:

$$r_i^{(t)} = \sigma \left(b_i^r + \sum_j U_{i,j}^r x_j^{(t)} + \sum_j W_{i,j}^r h_j^{(t)} \right)$$

Gated Recurrent Unit (GRU)

Reset and update gates can individually ignore parts of the state vector

- Update gates act like conditional leaky integrator that
- Update gates linearly integrate any dimension
 - choosing to copy it or completely
 - ignore it by replacing it with new target state value
 - target state is state the leaky integrator wants to converge towards
- Reset gates control which parts of the state get used to compute the next target state
- Reset gates introduce an additional nonlinear effect in the relationship between past and future state

Optimization for Long-Term Dependencies

When the parameter gradient is very large, gradient descent parameter update could throw parameters far away

A simple solution is **clipping gradients**

- Clip gradient from a minibatch element-wise before parameter update
- Clip the norm $\|\mathbf{g}\|$ of gradient \mathbf{g} before parameter update

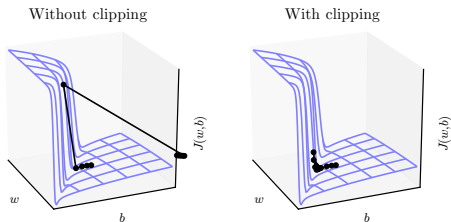


Figure 20: Gradient Clipping²².

²²| Goodfellow, Y Bengio, and A Courville. *Deep Learning*.

<http://www.deeplearningbook.org>. MIT Press, 2016.

Recurrent Neural Networks

Questions?

References



Felix A. Gers, Jürgen A. Schmidhuber, and Fred A. Cummins. “Learning to Forget: Continual Prediction with LSTM”. In: *Neural Comput.* 12.10 (Oct. 2000), pp. 2451–2471.



I Goodfellow, Y Bengio, and A Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.



Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Comput.* 9.8 (Nov. 1997), pp. 1735–1780.