

Computing the Gradient in an RNN

Sargur Srihari
srihari@buffalo.edu

Topics in Sequence Modeling

- Overview
 1. Unfolding Computational Graphs
 2. Recurrent Neural Networks
 3. Bidirectional RNNs
 4. Encoder-Decoder Sequence-to-Sequence Architectures
 5. Deep Recurrent Networks
 6. Recursive Neural Networks
 7. The Challenge of Long-Term Dependencies
 8. Echo-State Networks
 9. Leaky Units and Other Strategies for Multiple Time Scales
 10. LSTM and Other Gated RNNs
 11. Optimization for Long-Term Dependencies
 12. Explicit Memory

Topics in Recurrent Neural Networks

0. Overview
1. Teacher forcing and Networks with Output Recurrence
2. Computing the gradient in a RNN
3. RNNs as Directed Graphical Models
4. Modeling Sequences Conditioned on Context with RNNs

Training an RNN is straightforward

1. Determine gradient using backpropagation
 - Apply generalized backpropagation to the unrolled computational graph (algorithm repeated in next slide)
 - No specialized algorithms are necessary
 2. Gradient Descent
 - Gradients obtained by backpropagation are used with any general purpose gradient-based techniques to train an RNN
- Called BPTT: Back Propagation Through Time

General Backpropagation to compute gradient

To compute gradients \mathbb{T} of variable z wrt variables in computational graph \mathcal{G}

Algorithm 6.5 The outermost skeleton of the back-propagation algorithm. This portion does simple setup and cleanup work. Most of the important work happens in the `build_grad` subroutine of algorithm 6.6

Require: \mathbb{T} , the target set of variables whose gradients must be computed.

Require: \mathcal{G} , the computational graph

Require: z , the variable to be differentiated

Let \mathcal{G}' be \mathcal{G} pruned to contain only nodes that are ancestors of z and descendants of nodes in \mathbb{T} .

Initialize `grad_table`, a data structure associating tensors to their gradients

`grad_table[z] ← 1`

for \mathbf{V} in \mathbb{T} **do**

`build_grad(V, \mathcal{G} , \mathcal{G}' , grad_table)`

end for

Return `grad_table` restricted to \mathbb{T}

Build-grad function of Generalized Backprop

Algorithm 6.6 The inner loop subroutine `build_grad(V, G, G', grad_table)` of the back-propagation algorithm, called by the back-propagation algorithm defined in algorithm 6.5.

Require: \mathbf{V} , the variable whose gradient should be added to \mathcal{G} and `grad_table`.

Require: \mathcal{G} , the graph to modify.

Require: \mathcal{G}' , the restriction of \mathcal{G} to nodes that participate in the gradient.

Require: `grad_table`, a data structure mapping nodes to their gradients

if \mathbf{V} is in `grad_table` **then**

 Return `grad_table[V]`

end if

$i \leftarrow 1$

for \mathbf{C} in `get_consumers(V, G')` **do**

$\text{op} \leftarrow \text{get_operation}(\mathbf{C})$

$\mathbf{D} \leftarrow \text{build_grad}(\mathbf{C}, \mathcal{G}, \mathcal{G}', \text{grad_table})$

$\mathbf{G}^{(i)} \leftarrow \text{op.bprop}(\text{get_inputs}(\mathbf{C}, \mathcal{G}'), \mathbf{V}, \mathbf{D})$

$i \leftarrow i + 1$

end for

$\mathbf{G} \leftarrow \sum_i \mathbf{G}^{(i)}$

`grad_table[V] = G`

Insert \mathbf{G} and the operations creating it into \mathcal{G}

Return \mathbf{G}

Ob. `bprop(inputs, X, G)`

returns $\sum_i (\nabla_{\mathbf{x}} \text{op.f}(\text{inputs}))_i G_i$

Which is just an implementation of chain rule

If $\mathbf{y} = g(x)$ and $z = f(\mathbf{y})$ from chain rule $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

Equivalently, in vector notation where $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is the $n \times n$ Jacobian matrix of g

Example for how BPTT algorithm behaves

- BPTT: Back Propagation Through Time
- We provide an example of how to compute gradients by BPTT for the RNN equations

$$\begin{aligned}\mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}) \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)} \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)})\end{aligned}$$

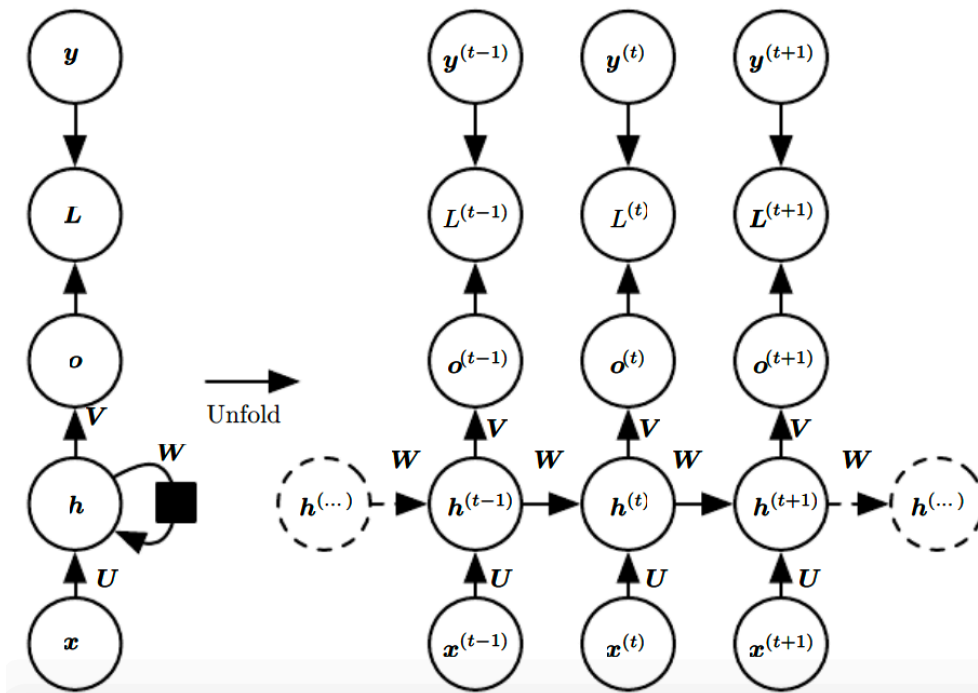
$$\begin{aligned}&L\left(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}\right) \\ &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}\right)\end{aligned}$$

Intuition of how BPTT algorithm behaves

- Ex: Gradients by BPTT for RNN eqns given above

Nodes of our computational graph include the parameters U, V, W, b, c as well as the sequence of nodes indexed by t for $\mathbf{x}^{(t)}, \mathbf{h}^{(t)}, \mathbf{o}^{(t)}$ and $L^{(t)}$

For each node \mathbf{N} we need to compute the gradient $\nabla_{\mathbf{N}} L$ recursively



Computing BPTT gradients recursively

- Start recursion with nodes immediately preceding final loss

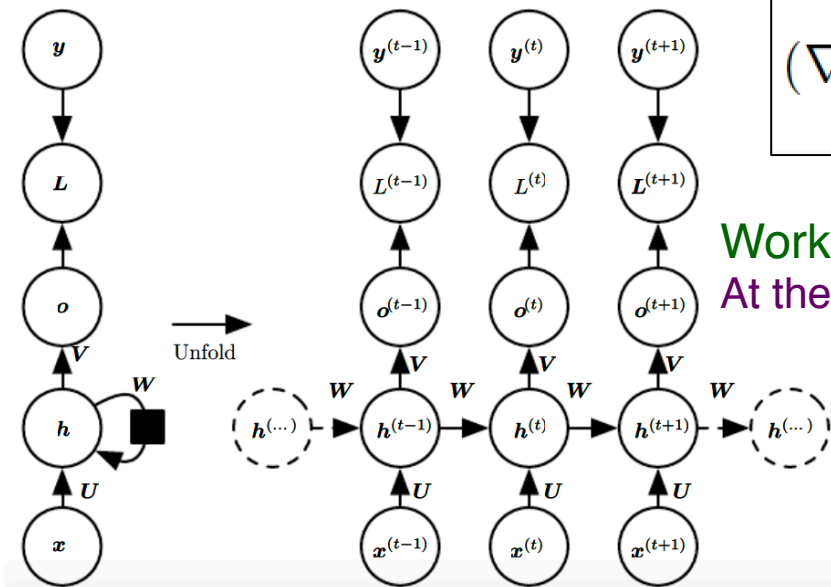
$$\frac{\partial L}{\partial L^{(t)}} = 1$$

- We assume that the outputs $o^{(t)}$ are used as arguments to softmax to obtain vector \hat{y} of probabilities over the output. Also that loss is the negative log-likelihood of the true target $y^{(t)}$ given the input so far.
- The gradient on the outputs at time step t is

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

Work backwards starting from the end of the sequence.
At the final time step τ , $h^{(\tau)}$ has only $o^{(\tau)}$ as a descendent. So

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



Gradients over hidden units

- Iterate backwards in time iterating to backpropagate gradients through time from $t = \tau - 1$ down to $t = 1$

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

- where $\text{diag} \left(1 - (\mathbf{h}^{(t+1)})^2 \right)$ indicates a diagonal matrix with elements within parentheses. This is the Jacobian of the hyperbolic tangent associated with the hidden unit i at time $t+1$

Gradients over parameters

- Once gradients on internal nodes of the computational graph are obtained, we can obtain gradients on the parameters
 - Because parameters are shared across time steps some care is needed

$$\begin{aligned} \nabla_{\mathbf{c}} L &= \sum_t \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L \\ \nabla_{\mathbf{b}} L &= \sum_t \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) \nabla_{\mathbf{h}^{(t)}} L \\ \nabla_{\mathbf{V}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{V} o_i^{(t)}} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top} \\ \nabla_{\mathbf{W}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)} h_i^{(t)}} \\ &= \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top} \\ \nabla_{\mathbf{U}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)} h_i^{(t)}} \\ &= \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top} \end{aligned}$$