

Recurrent Neural Networks

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

- Design patterns for RNNs
- RNN with recurrence between hidden units
- Forward propagation equations
- Loss function for a sequence
- RNN with recurrence from output units to hidden units

1. Teacher forcing for output-to-hidden RNNs

- Backward Propagation through time (BPTT)

2. Computing the gradient in a RNN

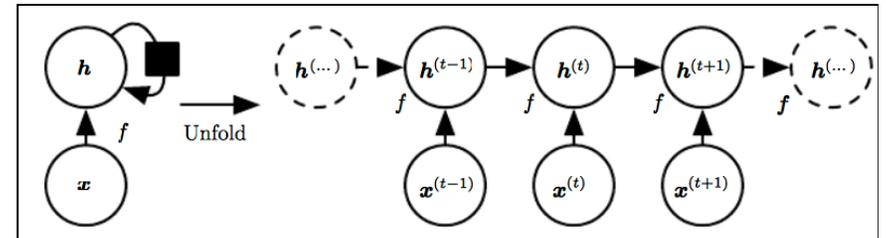
3. RNNs as Directed Graphical Models

4. Modeling Sequences Conditioned on Context with RNNs

Graph unrolling/parameter sharing in RNN design

- Process of graph unrolling:

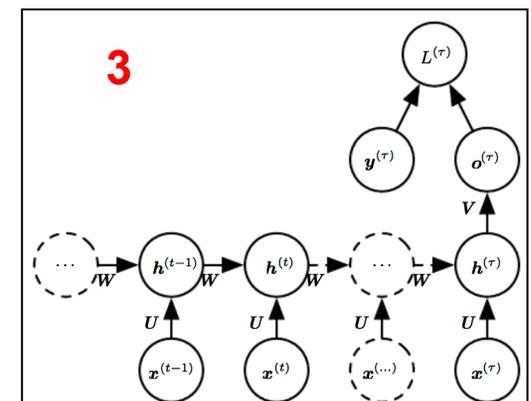
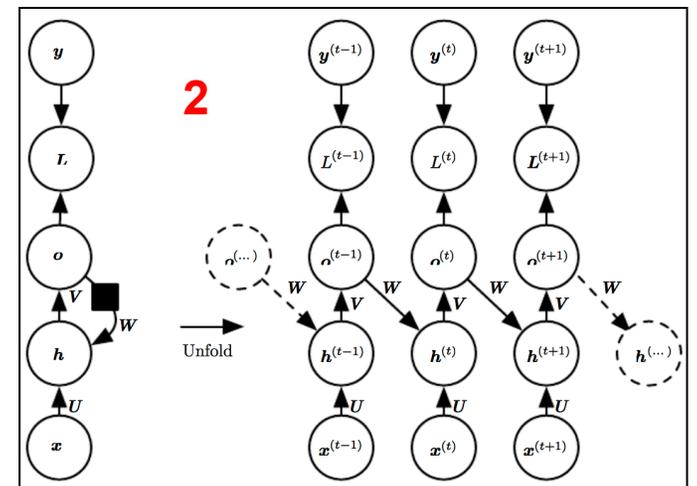
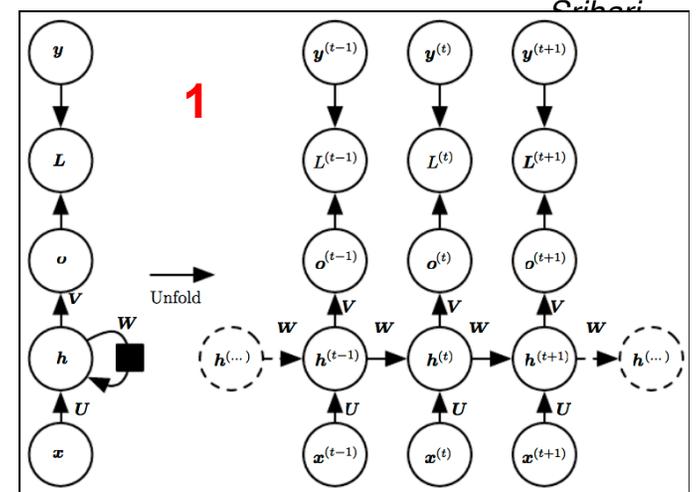
$$\begin{aligned} \mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}) \end{aligned}$$



- Function $g^{(t)}$ takes in whole past sequence $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$ as input and produces the current state but the unfolded recurrent structure allows us to factorize $g^{(t)}$ into repeated application of a function f
- It does not need a separate model $g^{(t)}$ for all possible time steps
- Process of unrolling and parameter sharing allows us to design a wide variety of recurrent neural networks
- Examples of important design patterns shown next

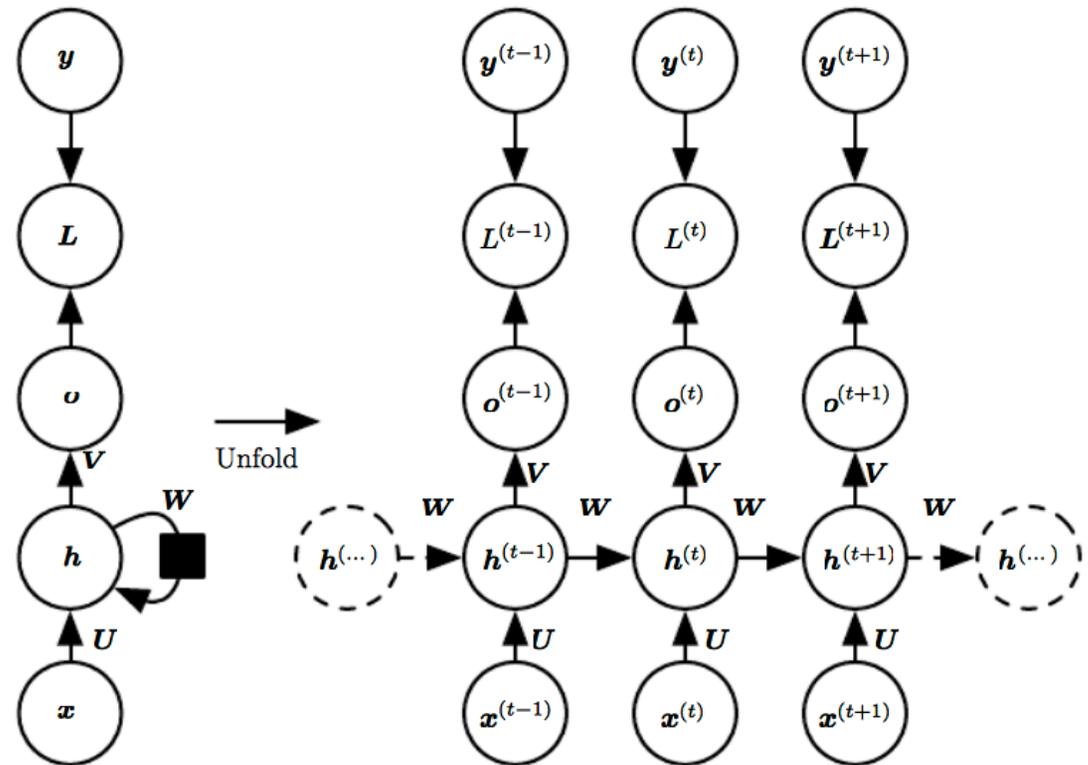
Three design patterns of RNNs

1. Produce output at each time step and have recurrent connections between hidden units
2. Produce output at each time step and have recurrent connections only from output at one time step to hidden units at next time step
3. Recurrent connections between hidden units to read entire input sequence and produce a single output
 - Can summarize sequence to produce a fixed size representation for further processing



Design 1: RNN with recurrence between hidden units

- Maps input sequence x to output o values
 - Loss L measures how far each o is from the corresponding target y
 - With softmax outputs we assume o is the unnormalized log probabilities
 - Loss L internally computes $y = \text{softmax}(o)$ and compares to target y



- Update equation

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

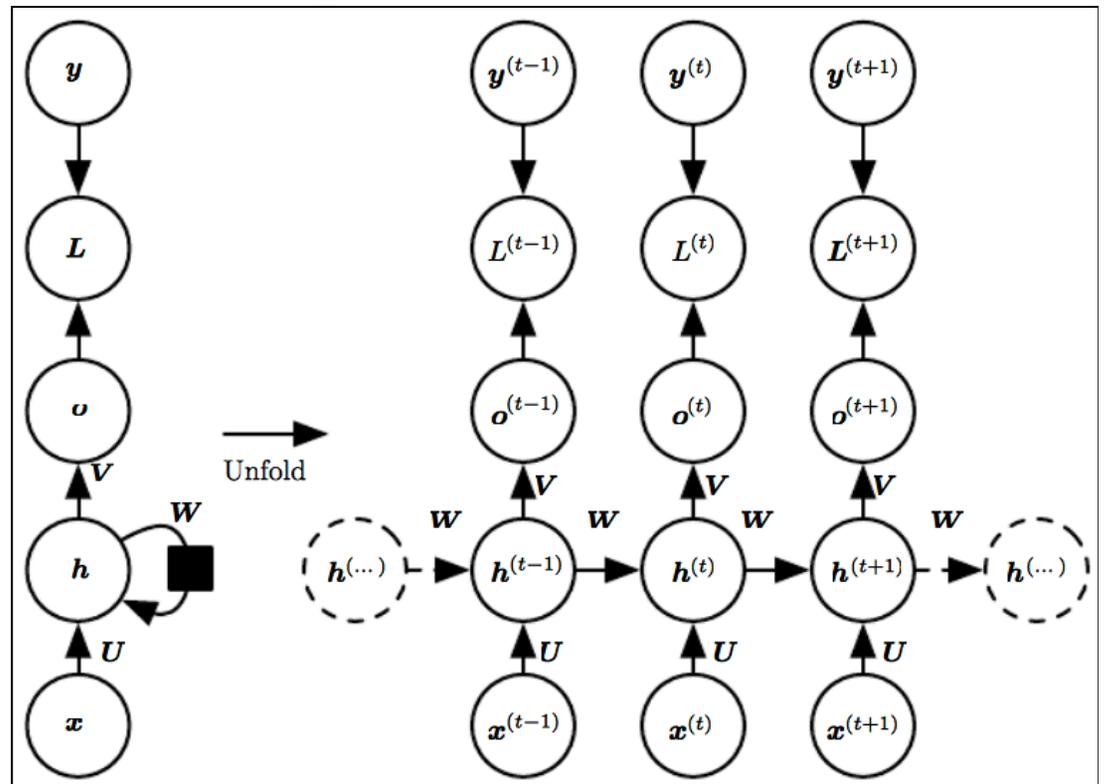
RNN with hidden unit recurrence is a Universal TM

- RNN with hidden unit connections together with

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

is Universal, i.e.,

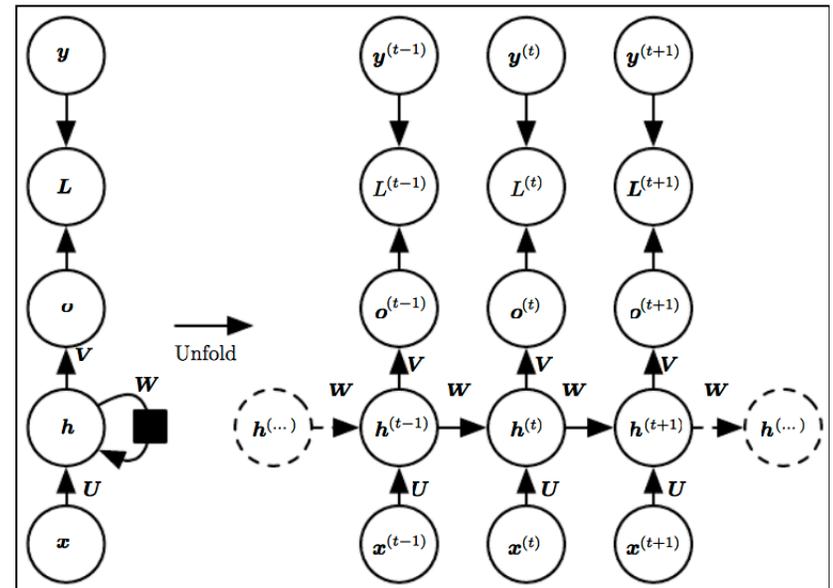
- Any function computable by a TM is computable by such an RNN of finite size



- Output can be read from the RNN after no. of steps asymptotically linear in no. of steps used by the TM and asymptotically linear in length of input

Forward prop for RNN with hidden unit recurrence

- This design does not specify
 - activation functions for hidden units,
 - form of output and
 - loss function
- Assume for this graph:
 - Hyperbolic tangent activation function
 - Output is discrete
 - E.g., RNN predicts words/characters
 - Natural way to represent discrete vars:
 - Output o gives unnormalized log probabilities \hat{y} of each possible value
 - Can apply `softmax` as a postprocessing step to obtain vector of normalized probabilities over the output
- Forward prop begins with specification of initial state $h^{(0)}$



Forward Prop Equations for RNN with hidden recurrence

- Begins with initial specification of $\mathbf{h}^{(0)}$
- Then for each time step from $t=1$ to $t=\tau$ we apply the following update equations

$$\begin{array}{l} \mathbf{o}^{(t)} = \mathbf{c} + V\mathbf{h}^{(t)} \\ \mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)}) \\ \mathbf{a}^{(t)} = \mathbf{b} + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} \end{array}$$

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

- Where the parameters are:
 - bias vectors \mathbf{b} and \mathbf{c}
 - weight matrices U (input-to-hidden), V (hidden-to-output) and W (hidden-to-hidden) connections

Loss function for a given sequence

- This is an example of an RNN that maps an input sequence to an output sequence of the same length
- Total loss for a given sequence for a given sequence of \mathbf{x} values with a sequence of \mathbf{y} values is the sum of the losses over the time steps
- If $L^{(t)}$ is the negative log-likelihood of $\mathbf{y}^{(t)}$ given $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ then

$$L\left(\left\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\right\}, \left\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t)}\right\}\right) = \sum_t L^{(t)}$$

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

- where p_{model} is given by reading the entry for $\mathbf{y}^{(t)}$ from the model's output vector $\hat{\mathbf{y}}^{(t)}$

Computing Gradient of Loss Function

- Computing gradient of this loss function wrt parameters is expensive
 - It involves performing a forward propagation pass moving left to right through the unrolled graph followed by a backward propagation moving right to left through the graph
- Run time is $O(\tau)$ and cannot be reduce by parallelization
 - Because forward propagation graph is inherently sequential
 - Each time step computable only after previous step
 - States computed during forward pass must be stored until reused in the backward pass
 - So memory cost is also $O(\tau)$
- Backpropagation applied to the unrolled graph with $O(\tau)$ cost is called *Backward Propagation through time* (BPTT)

Backward Propagation through Time (BPTT)

- RNN with hidden unit recurrence is very powerful but also expensive to train
- Is there an alternative?