

RNNs: Teacher Forcing

Sargur Srihari
srihari@buffalo.edu

Topics in Sequence Modeling

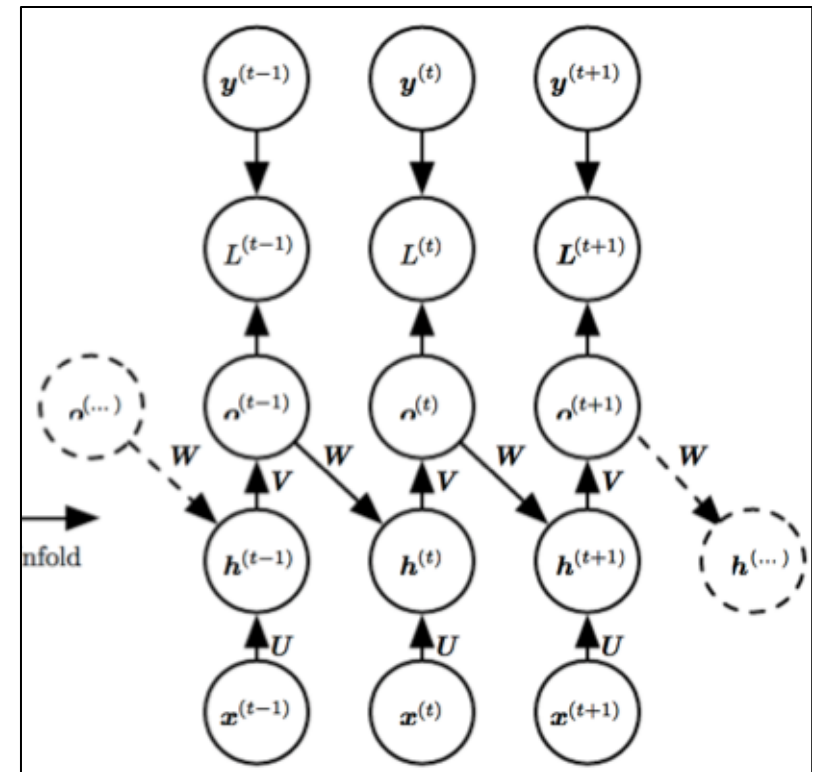
- Recurrent Neural Networks
 1. Unfolding Computational Graphs
 2. Recurrent Neural Networks
 3. Bidirectional RNNs
 4. Encoder-Decoder Sequence-to-Sequence Architectures
 5. Deep Recurrent Networks
 6. Recursive Neural Networks
 7. The Challenge of Long-Term Dependencies
 8. Echo-State Networks
 9. 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

- 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 and Networks with Output Recurrences
- 2. Computing the gradient in a RNN
- 3. Recurrent Networks as Graphical Models
- 4. Modeling Sequences Conditioned on Context with RNNs

Models with recurrence from output to hidden

- Less powerful than with hidden-to-hidden recurrent connections
 - It cannot simulate a universal TM
 - It requires that the output capture all information of past to predict future
- Advantage
 - In comparing loss function to output all time steps are decoupled
 - Each step can be trained in isolation
 - Training can be parallelized
 - Gradient for each step t computed in isolation
 - No need to compute output for the previous step first, because training set provides ideal value of output
- Can be trained with Teacher Forcing
 - Described next



Teacher Forcing

- Used for training models with recurrent connections from their outputs
- Teacher forcing during training means
 - Instead of summing activations from incoming units (possibly erroneous)
 - Each unit sums correct teacher activations as input for the next iteration

Teacher Forcing Procedure

- Teacher forcing is a procedure for training RNNs with output-to-hidden recurrence
- It emerges from the maximum likelihood criterion
- During training time model receives ground truth output $y^{(t)}$ as input at time $t+1$
 - We can see this by examining a sequence with two time steps (next)

Maximum Likelihood in Teacher Forcing

- Consider a sequence with two time steps
- The conditional maximum likelihood criterion is

$$\log (p(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)})) = \log p(\mathbf{y}^{(2)} | \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) + \log p(\mathbf{y}^{(1)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)})$$

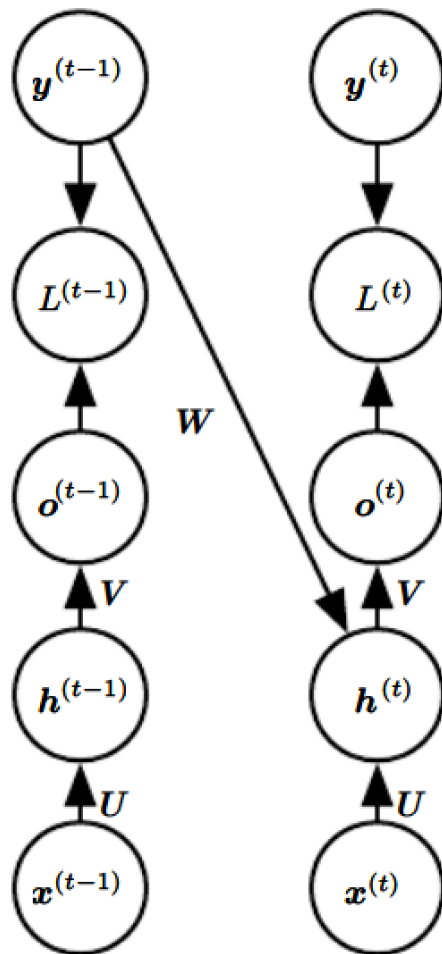
- At time $t=2$, model is trained to maximize conditional probability of $\mathbf{y}^{(2)}$ given both the \mathbf{x} sequence so far and the previous \mathbf{y} value from the training set.
 - We are using $\mathbf{y}^{(1)}$ as teacher forcing, rather than only $\mathbf{x}^{(i)}$.
- Maximum likelihood specifies that during training, rather than feeding the model's own output back to itself, target values should specify what the correct output should be

Illustration of Teacher Forcing

- Teacher Forcing is a training technique applicable to RNNs that have connections from output to hidden states at next time step

Train time:

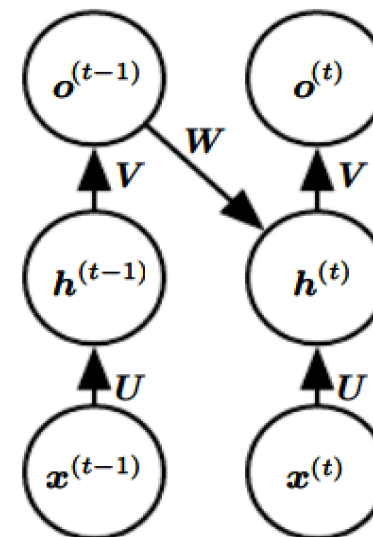
We feed the correct output $y^{(t)}$ (from teacher) drawn from the training set as input to $h^{(t+1)}$



Test time:

True output is not known.

We approximate the correct output $y^{(t)}$ with the model's output $o^{(t)}$ and feed the output back to the model

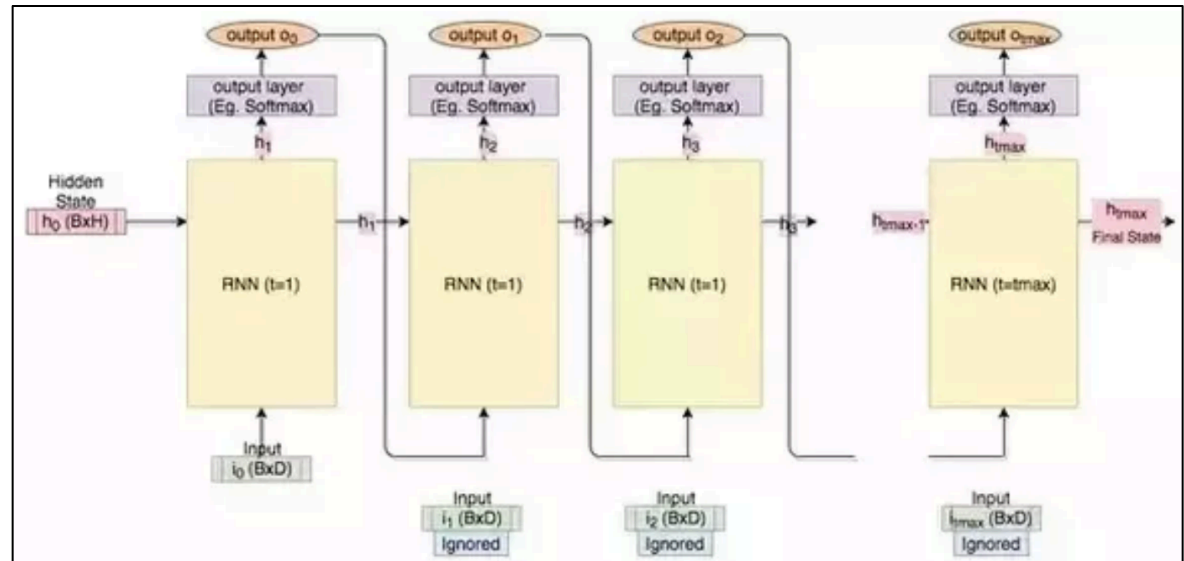


RNNs without and with teacher forcing

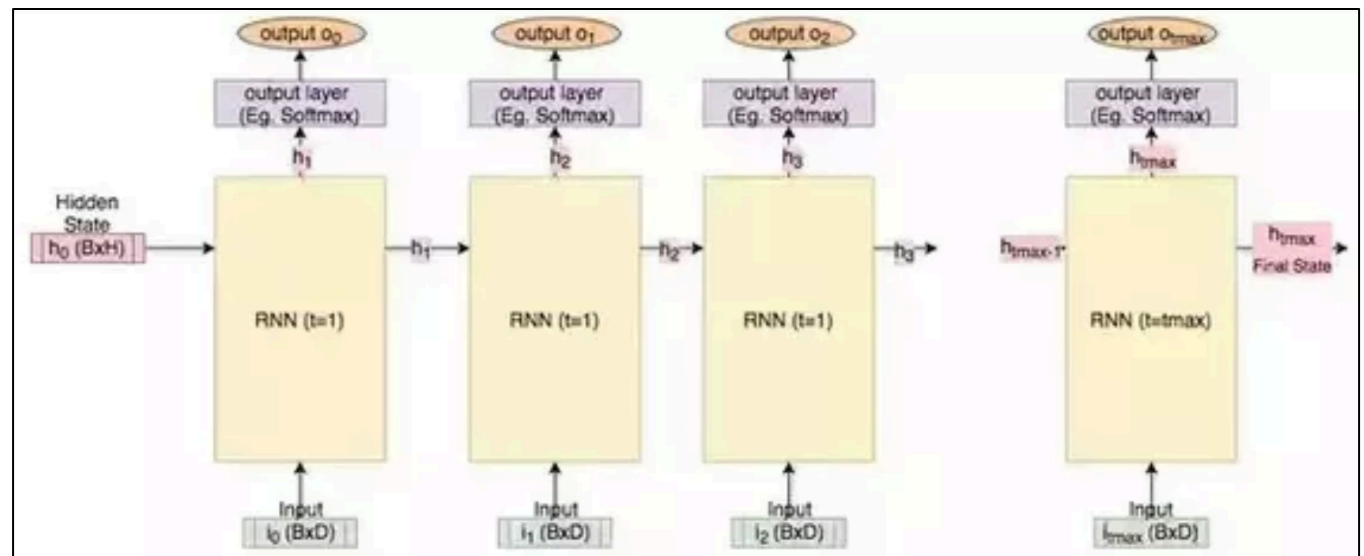
Without Teacher forcing:

Output of previous time step acts as input during RNN training

RNN learns to generate next output in a sequence based upon the previous value



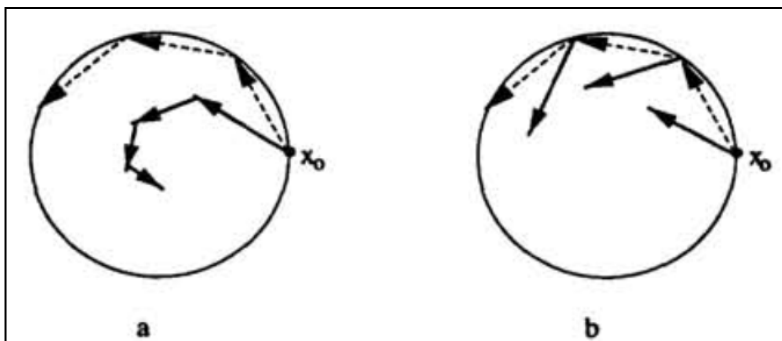
With Teacher forcing:



Source: <https://www.quora.com/What-is-the-teacher-forcing-in-RNN>

Visualizing Teacher Forcing

- Imagine that the network is learning to follow a trajectory
- It goes astray (because the weights are wrong) but teacher forcing puts the net back on its trajectory
 - By setting the state of all the units to that of teacher's.



(a) Without teacher forcing, trajectory runs astray (solid lines) while the correct trajectory are the dotted lines

(b) With teacher forcing trajectory corrected at each step

Teacher forcing for hidden-to-hidden

- We originally motivated teacher forcing as allowing us to avoid backpropagation through time in models that lack hidden-to-hidden connections
- Teacher forcing can still be applied to models that have hidden-to-hidden connections
 - As long as they have connections from the output at one time step to values computed at the next time step
 - As soon as the hidden units become functions of earlier time steps, BPTT algorithm is necessary

Training with both Teacher Forcing and BPTT

- Some models may be trained with both Teacher forcing and Backward Propagation through time (BPTT)
 - When there are both hidden-to-hidden recurrences as well as output-to-hidden recurrences

Disadvantage of Teacher Forcing

- If network is to be used in an *open-loop* mode with network outputs (or samples from the output distribution) fed back as inputs
 - In this case the kind of inputs that it will see during training time could be quite different from that it will see at test time
- Solutions to mitigate this problem:
 1. Train with both teacher-forced inputs and free running inputs
 - E.g., predicting the correct target a no of steps in the future through the unfolded recurrent output-to-input paths
 - Thus network can learn to take into account input conditions not seen during training
 2. Mitigate the gap between inputs seen at training time and test time by *generating* values as input
 - This approach exploits a *curriculum learning* strategy to gradually use more of the generated values as input

Professor Forcing

- Teacher Forcing trains RNNs by supplying observed sequence values as inputs during training and using the network's own one-step-ahead predictions to do multi-step sampling
- Professor Forcing algorithm
 - Uses adversarial domain adaptation to encourage the dynamics of the RNN to be the same when training the network and when sampling from the network over multiple time steps.
 - Applied to
 - language modeling
 - vocal synthesis on raw waveforms
 - handwriting generation
 - image generation
- Professor Forcing is an adversarial method for learning generative models
 - It is closely related to Generative Adversarial Networks