

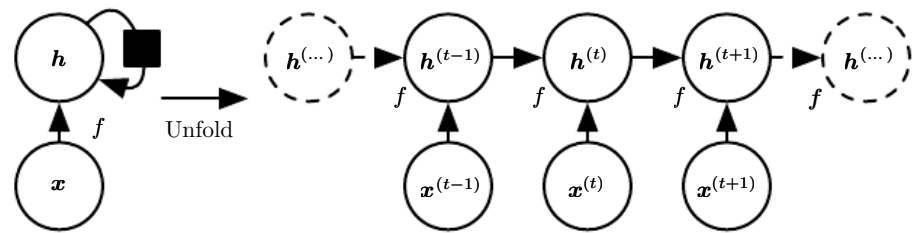
Recurrent Neural Network

CHEN Si

Contents

1	Definition	1
1.1	Motivation	2
1.2	Advantages	2
1.3	Disadvantages	3
1.4	How to determine the length τ	3
2	Taxonomy	3
2.1	3
2.2	5
2.3	6
2.4	8
2.5	9
3	Back Propagation Through Time (BPTT)	9

1 Definition



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

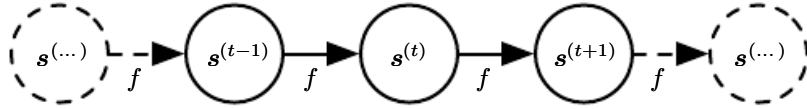
- Hidden units \mathbf{h} represents
 1. the state of the dynamical system.

2. lossy summary of the task-relevant aspects of the past sequence of inputs up to t , when the recurrent network is trained to perform a task that requires predicting the future from the past.
- Black square indicates that an interaction takes place with a delay of a single time step.

1.1 Motivation

Dynamical System

Unfolded Computational Graph



where $s^{(t)}$ is called the **state** of the system.

Recurrent Definition

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

Unfolded Definition (up to step 3)

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

Dynamical System with External Signal $x^{(t)}$

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

1.2 Advantages

1. Learning just a single model f :
 - The learned model always has the same input size, regardless of the sequence length.
 - Parameter sharing: It is possible to use the same transition function f with the same parameters at every time step.
 - Allows generalization to sequence lengths that does not appear in the training set.
 - Requires much fewer training examples to estimate.

1.3 Disadvantages

1. Optimizing the parameters may be difficult, because of the reduced number of parameters.

1.4 How to determine the length τ

1. Add a special symbol corresponding to the end of a sequence.
2. Introduce an extra Bernoulli output to the model that represents the decision to either continue generation or halt generation at each time step. (a more general way)
3. Add an extra output to the model that predicts the integer τ itself.

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

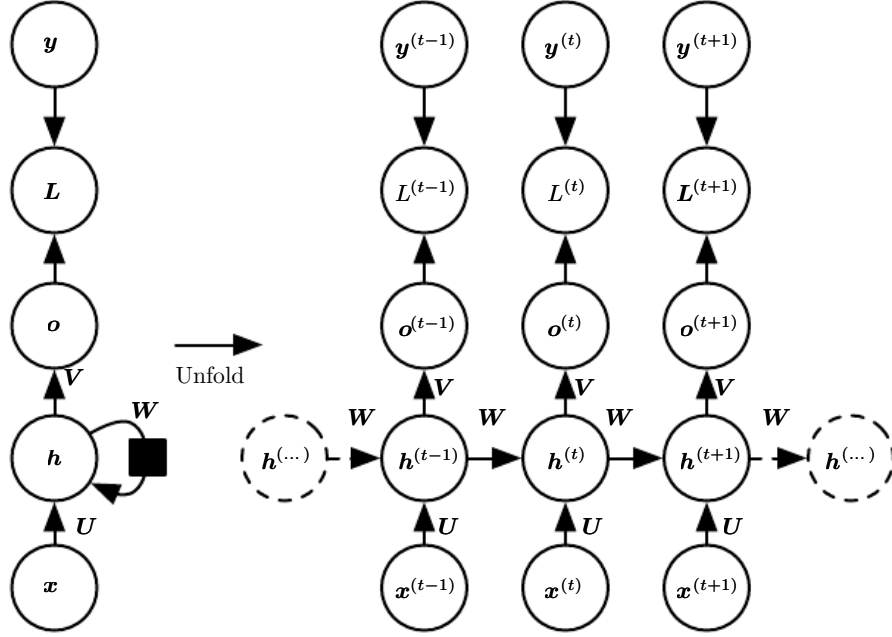
2 Taxonomy

2.1

Recurrent networks that produce an output at each time step and have recurrent connections between hidden units.

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

with a initial state $\mathbf{h}^{(0)}$.



Characteristics

- Maps an input sequence to an output sequence of the same length.
- Loss:

$$\begin{aligned}
 & L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}) \\
 &= \sum_t L^{(t)} \\
 &= - \sum_t \log p_{\text{model}}(\mathbf{y}^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\})
 \end{aligned} \tag{2}$$

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

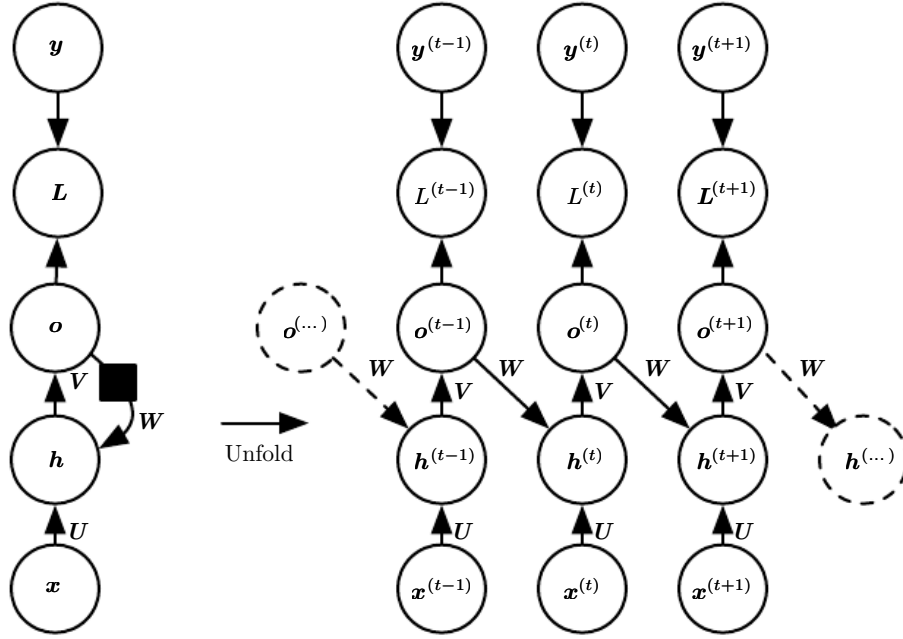
- Conditional independence assumption: The conditional distribution

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

- Runtime is $O(\tau)$, memory cost is $O(\tau)$.
- Only able to represent distributions in which the \mathbf{y} values are conditionally independent from each other given the \mathbf{x} values.

2.2

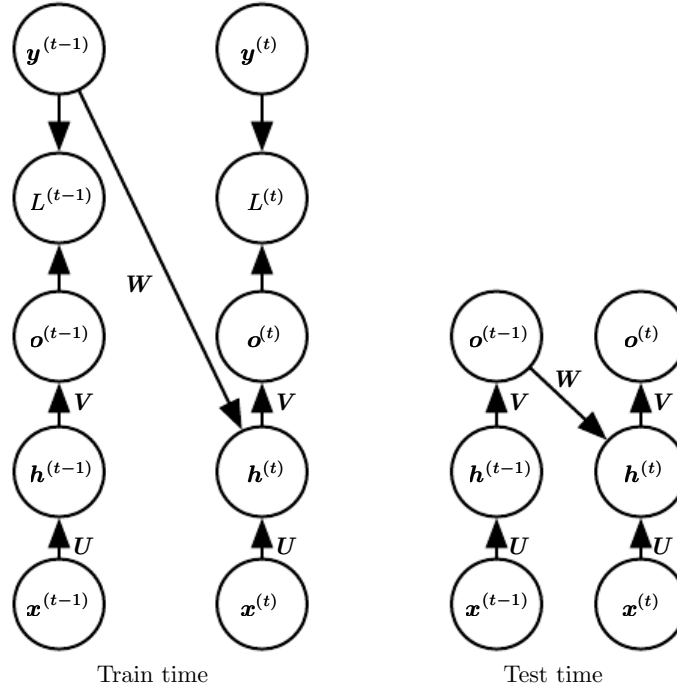
Recurrent networks that produce an output at each time step and have recurrent connections only from the output at one time step to the hidden units at the next time step.



Characteristics

- Strictly less powerful because it lacks hidden-to-hidden recurrent connections. (e.g. it cannot simulate a universal Turing machine.)
- Training can be **parallelized** because all the time steps are decoupled (the training set provides the ideal value of previous output).
- Avoids back-propagation through time.
- Teacher forcing: For example, for a sequence with two time steps

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



Disadvantages

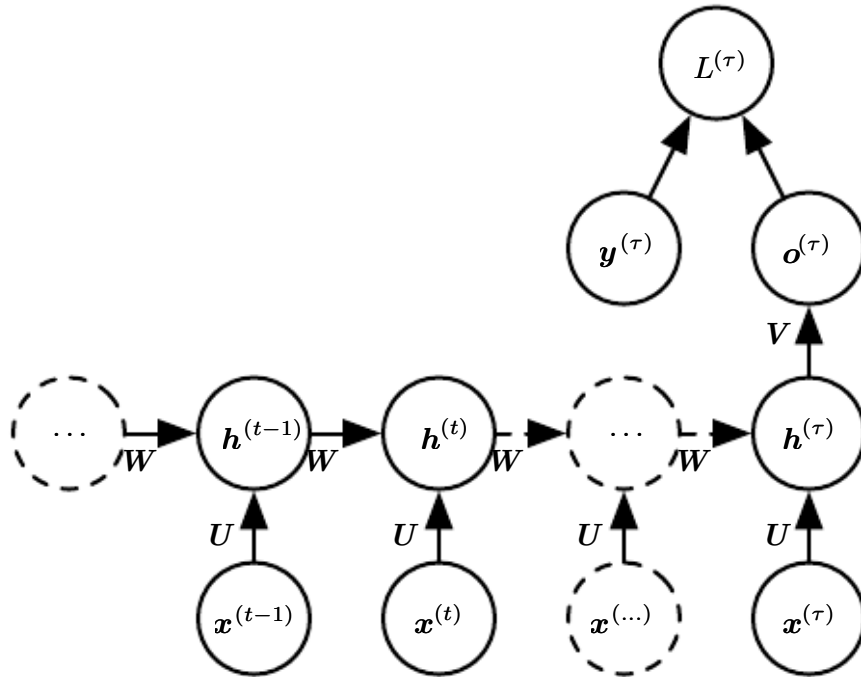
1. The disadvantages of strict teacher forcing (no BPTT) is that the inputs during training could be quite different from the inputs during testing (**closed-loop** mode).

Solutions:

- (a) Train with both teacher-forced inputs and free-running inputs
- (b) Randomly chooses to use generated values or actual data values as input (curriculum learning strategy: gradually use more of the generated values as input).

2.3

Recurrent networks with recurrent connections between hidden units, that read an entire sequence and then produce a single output.



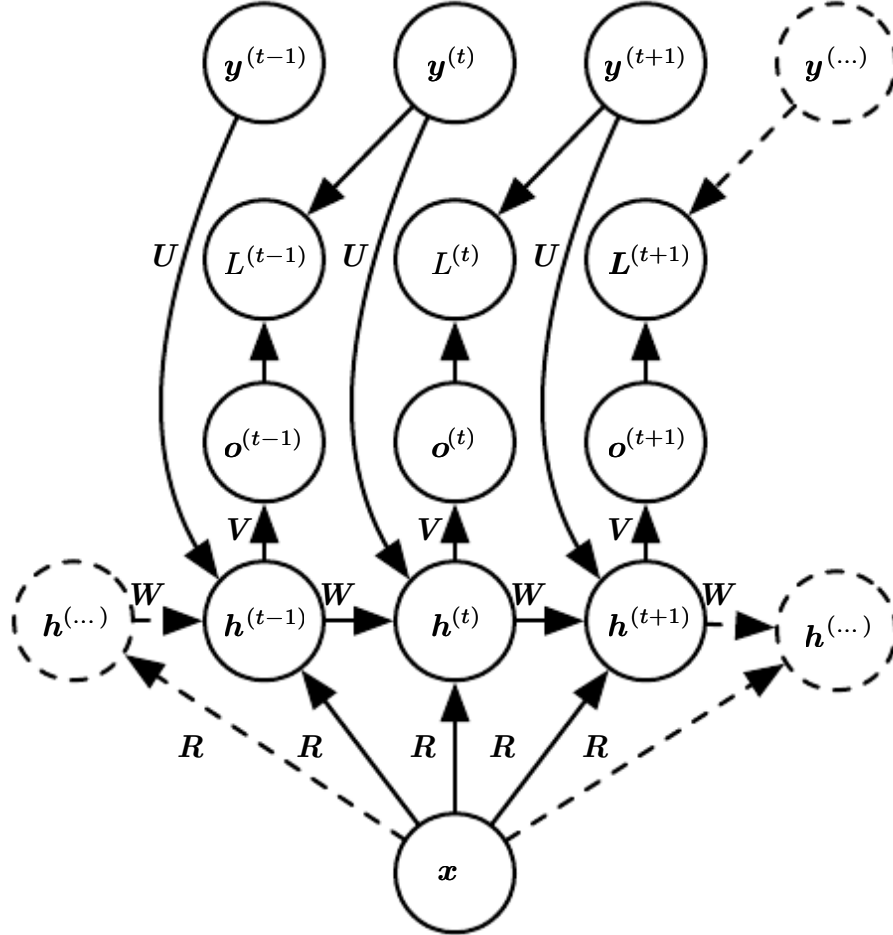
Characteristics

- There might be a target right at the end, or the gradient on the output $o^{(t)}$ can be obtained by back-propagating from further downstream modules.

Applications

1. Summarize a sequence and produce a fixed-size representation used as input for further processing.

2.4



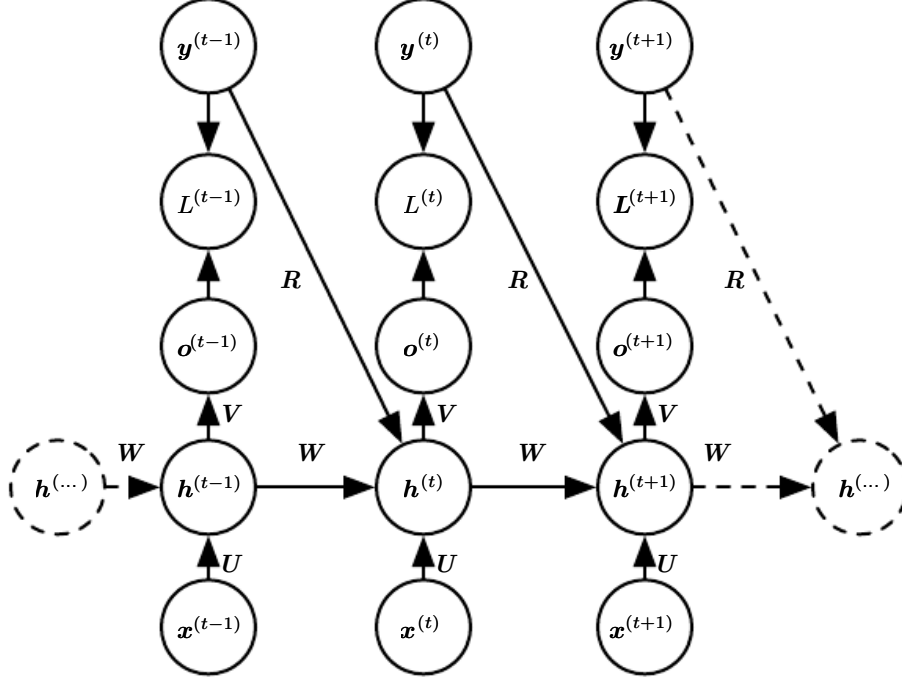
Characteristics

- $\mathbf{x}^\top \mathbf{R}$ is effectively a new bias parameter.
- Each element $\mathbf{y}^{(t)}$ of the observed output sequence serves both as input (for the current time step) and, during training, as target (for the previous time step).

Applications

1. Image captioning, where a single image is used as input to a model that then produces a sequence of words describing the image.

2.5



3 Back Propagation Through Time (BPTT)

Here we talk about the BPTT of equation 1 as forward propagation and equation 2 as loss function. (Note: here $\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$)

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

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

$$\forall t = 1, \dots, \tau - 1$$

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

4.

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