

Recurrent Neural Networks

Seminar Principles of Data Mining and Learning Algorithms

Gular Shukurova
Sheikh Mastura Farzana

Outline

- Unfolding Computational Graphs (10.1)
- Recurrent Neural Networks (10.2.1)
- Gradient Computation (10.2.2)
- RNNs as Directed Graphical Models (10.2.3)
- Modeling Sequences Conditioned on Context (10.2.4)
- Long Term Dependencies (10.7)
- Gated RNNs (10.10)

Unfolding Computational Graphs

- Expressing a recurrent computation into a computational graph

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

For $t=3$

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

Example

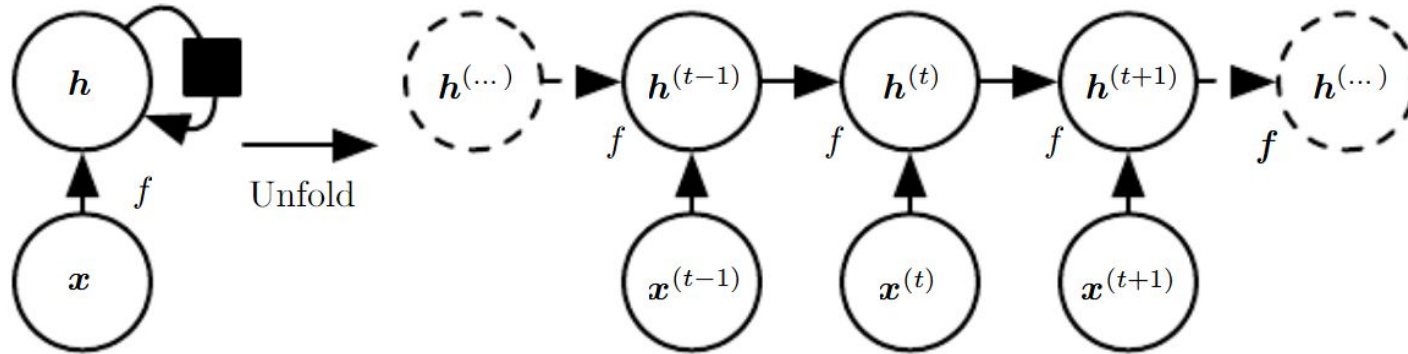


Figure: (left) Circuit Diagram, (right) unfolded computational graph, each node associated to a single timestep.

Computation

Rewriting the equation from previous slide with $h(t)$:

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

The unfolded recurrence after t steps represented with a function $g(t)$:

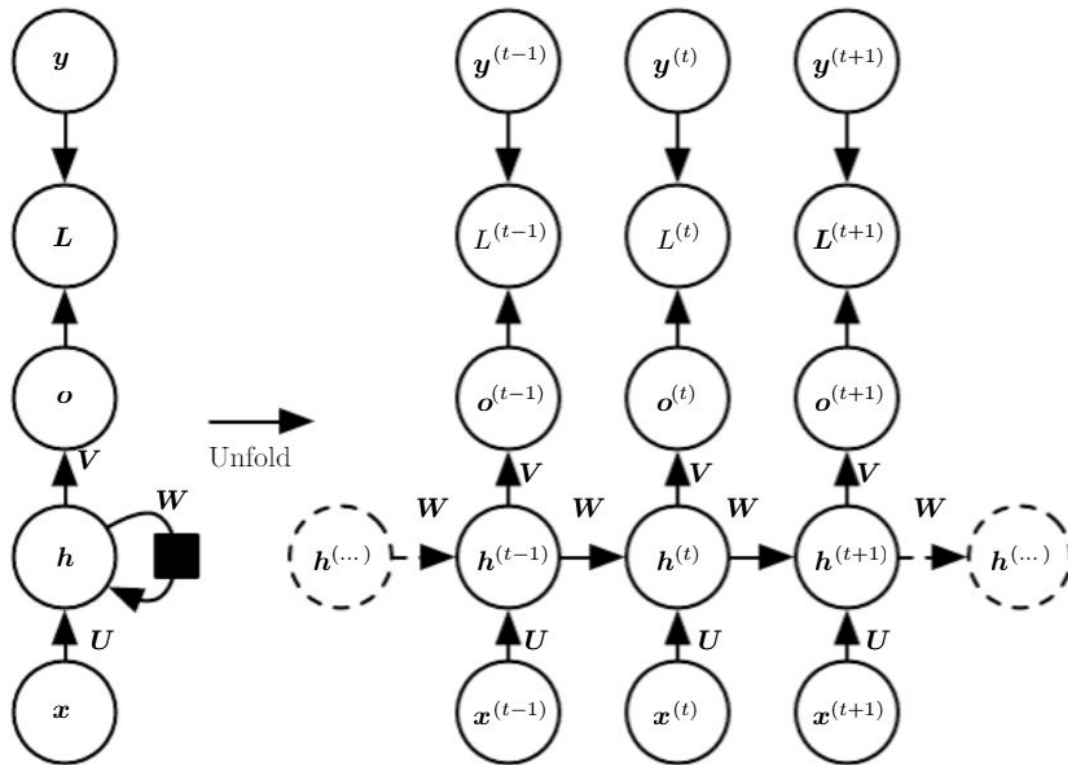
$$\begin{aligned}\mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}).\end{aligned}$$

Recurrent Neural Networks

- **Recurrent Neural Networks (RNNs)** are a class of neural networks for processing sequential data.
- RNNs use **feedback loops** to process a sequence of data that allow information to persist.
- Reducing the complexity of parameters by **parameter sharing**
- A powerful tool in applications like text processing, speech recognition, language translation and DNA sequences, where the output depends on the previous computations.

RNNs by examples

- Example #1



from [1]

RNNs by examples (cont)

- Forward propagation, $t \in [1, \tau]$:

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

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

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

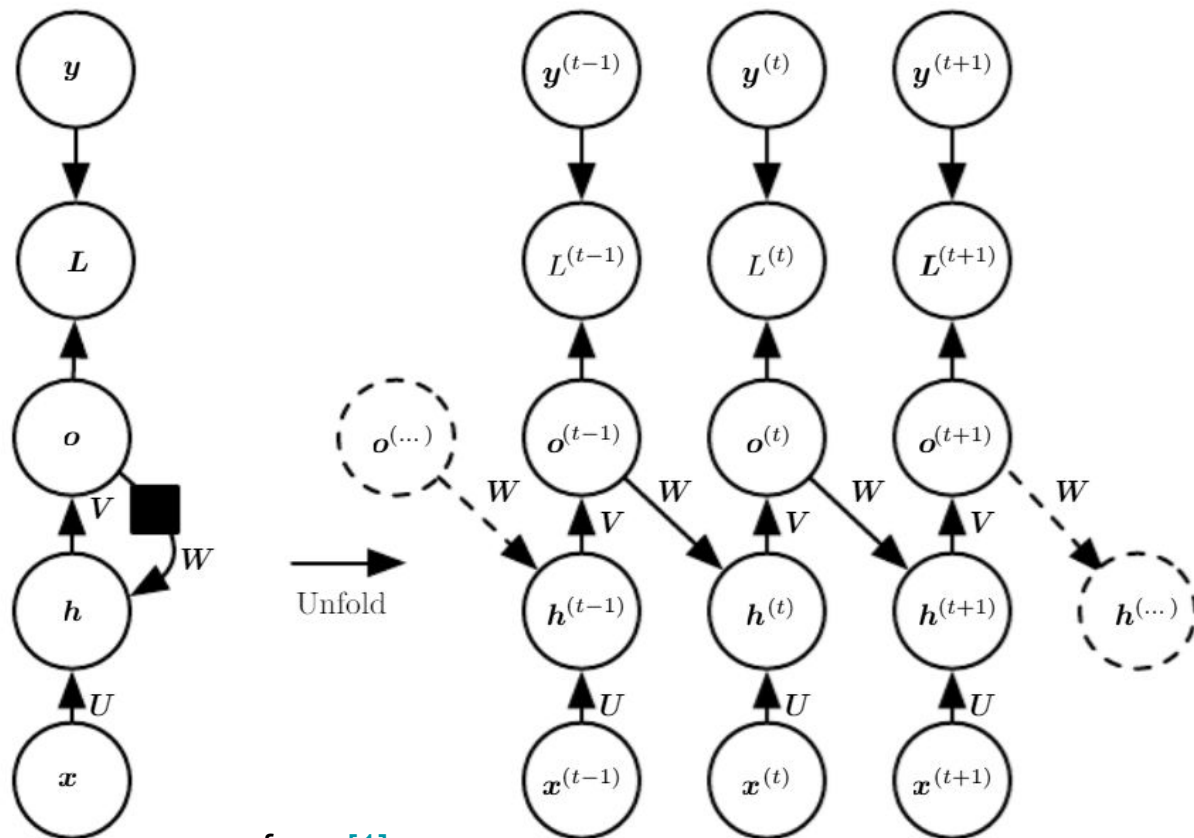
$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}),$$

- Total loss:

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

RNNs by examples (cont)

- Example #2



from [1]

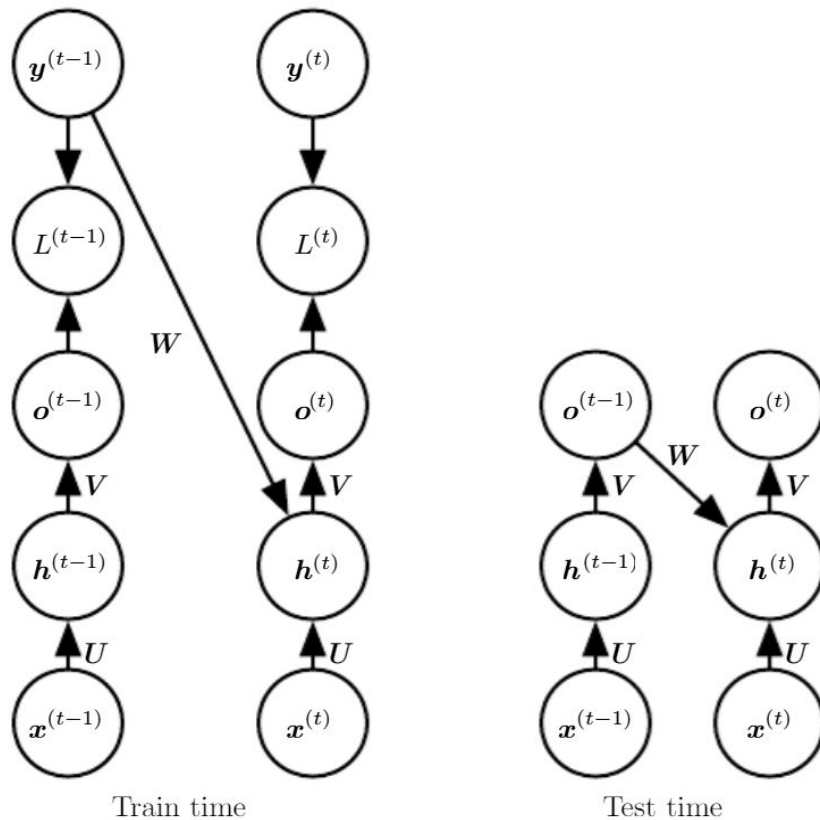
Teacher Forcing

- Conditional maximum likelihood criterion:

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

- Advantage: to avoid **BPTT** in models that lack hidden-to-hidden connections
- Disadvantage: works poorly in **open-loop** mode
 - 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

Teacher Forcing (cont)



from [1]

Computing the gradient

- Based on equations on [slide 5](#)

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

- $t = \tau$:

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

- $t \in [\tau-1, 1]$

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

Computing the gradient (cont)

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

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

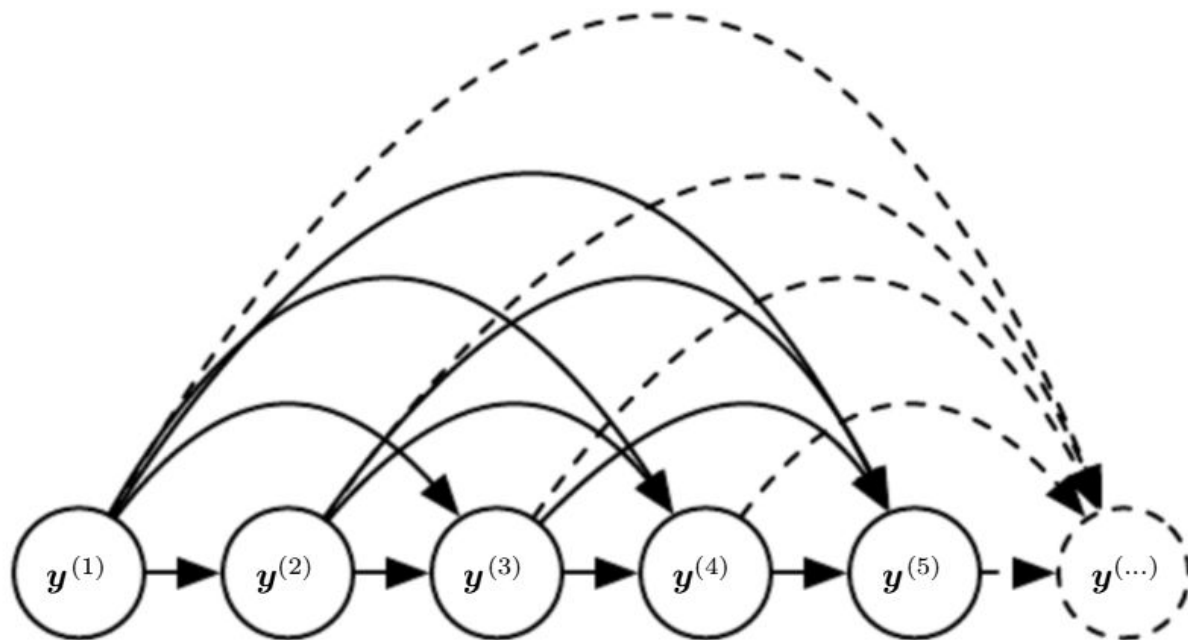
$$\nabla_{\mathbf{v}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{v}^{(t)}} o_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top},$$

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

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

RNN as Directed Graphical Models

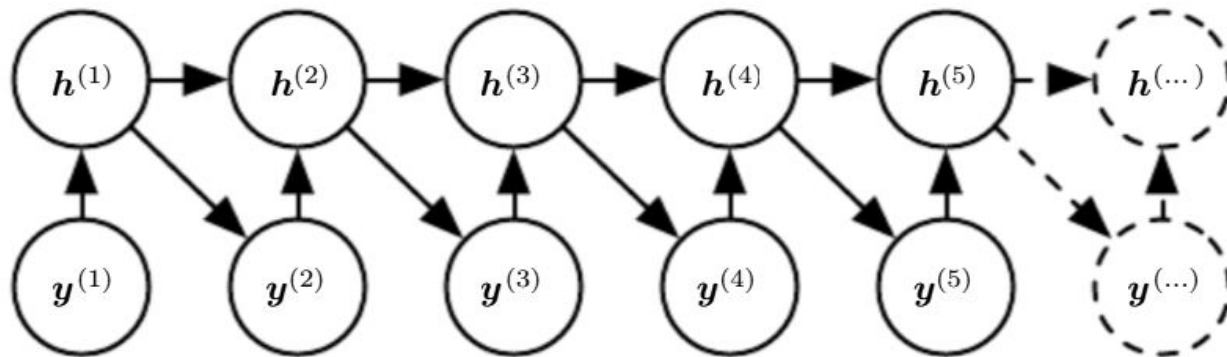
- Ignoring the hidden units
- inefficient



from [\[1\]](#)

RNN as Directed Graphical Models (cont)

- very efficient parametrization

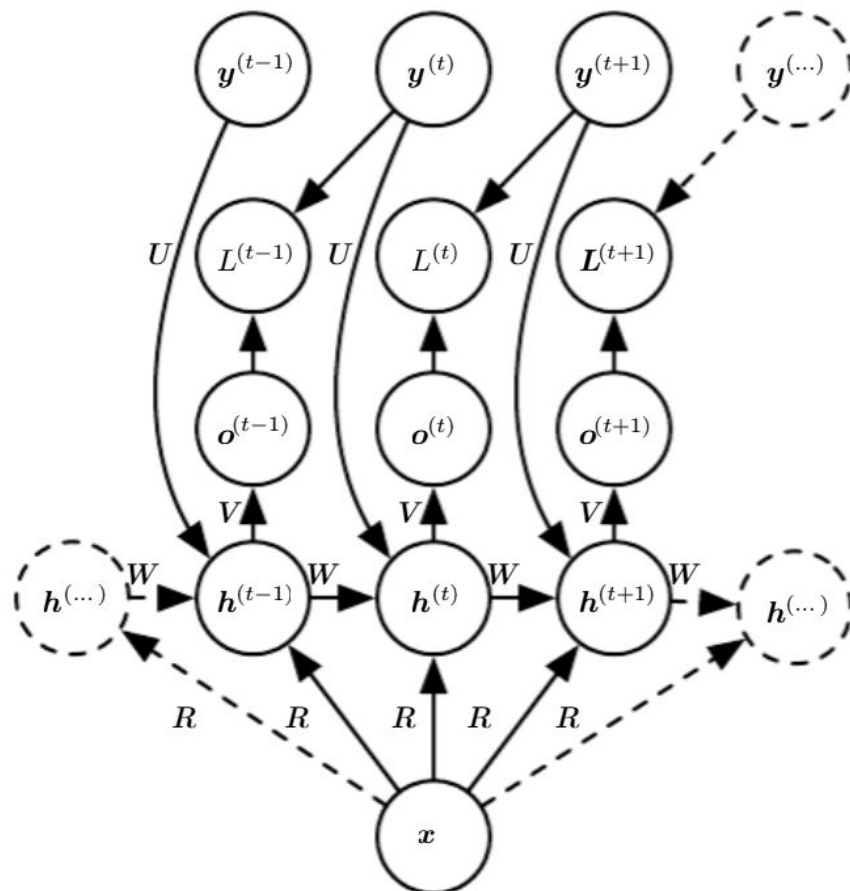


RNN as Directed Graphical Models (cont)

- Determining the length of the sequence
 - Special symbol at the end of the sequence
 - extra Bernoulli output
 - Predicting sequence length τ

Modeling Sequences Conditioned on Context

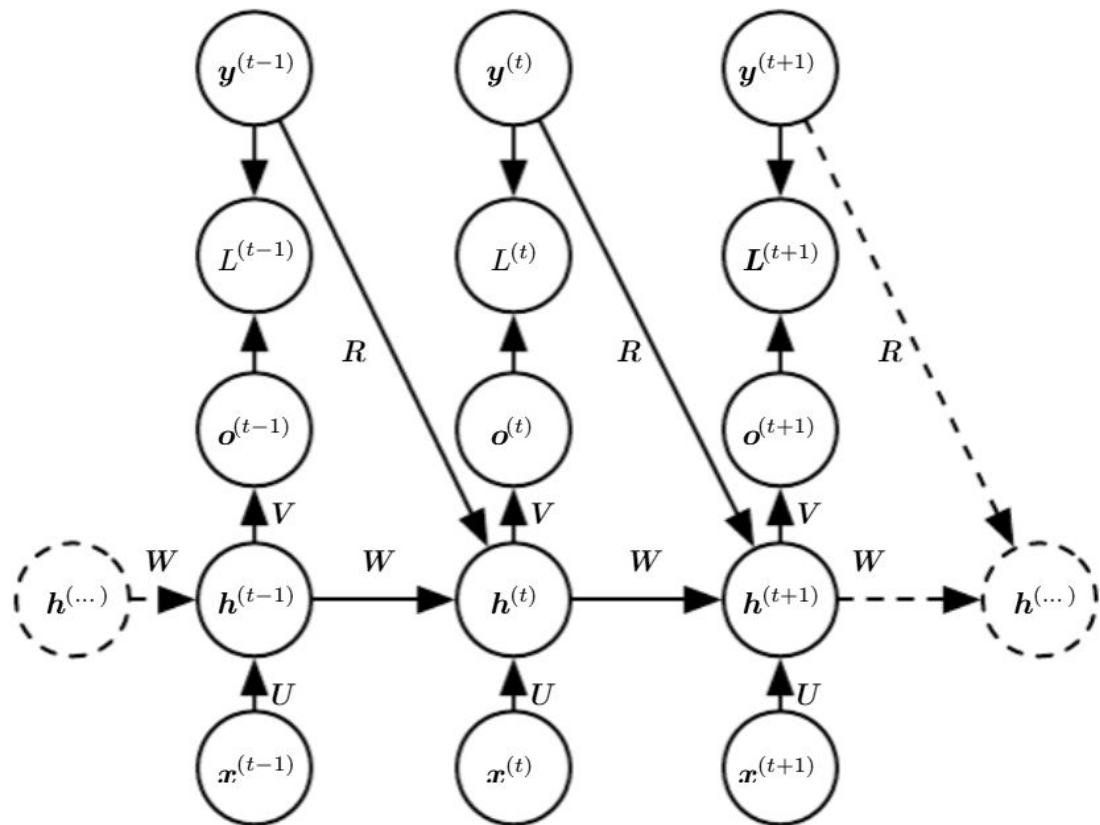
- A single vector as input



from [1]

Modeling Sequences Conditioned on Context (cont)

- A sequence of vectors as input

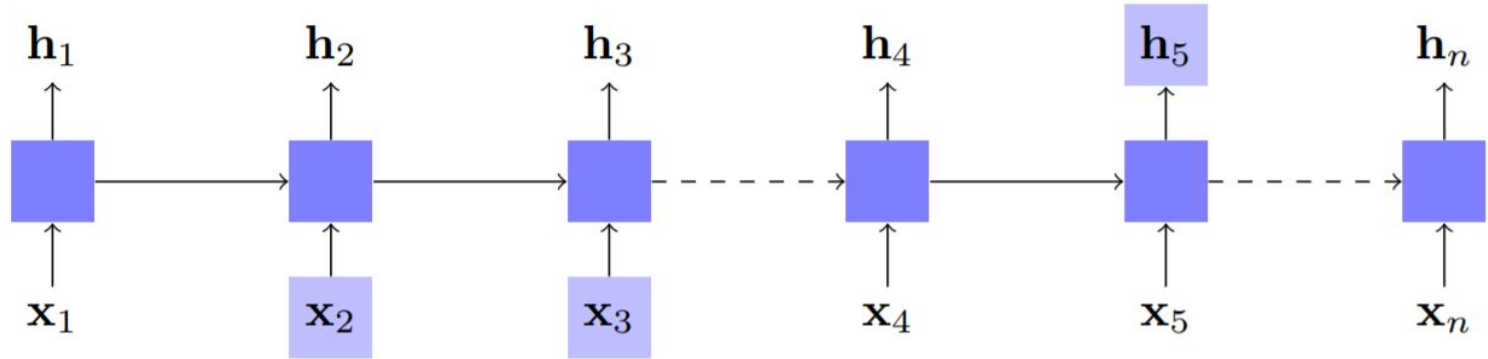


Long Term Dependencies

- Vanishing and exploding gradients in long-term propagation.
- Exponentially smaller magnitude of gradient for long term dependencies.
- Gradient based optimization is difficult.

Why do we want longterm dependencies?

I grew up in France ... I speak fluent French ...



Reference: RNNII, DLVR, Lecture by Dr. Michael Weinmann, Informatik, University of Bonn

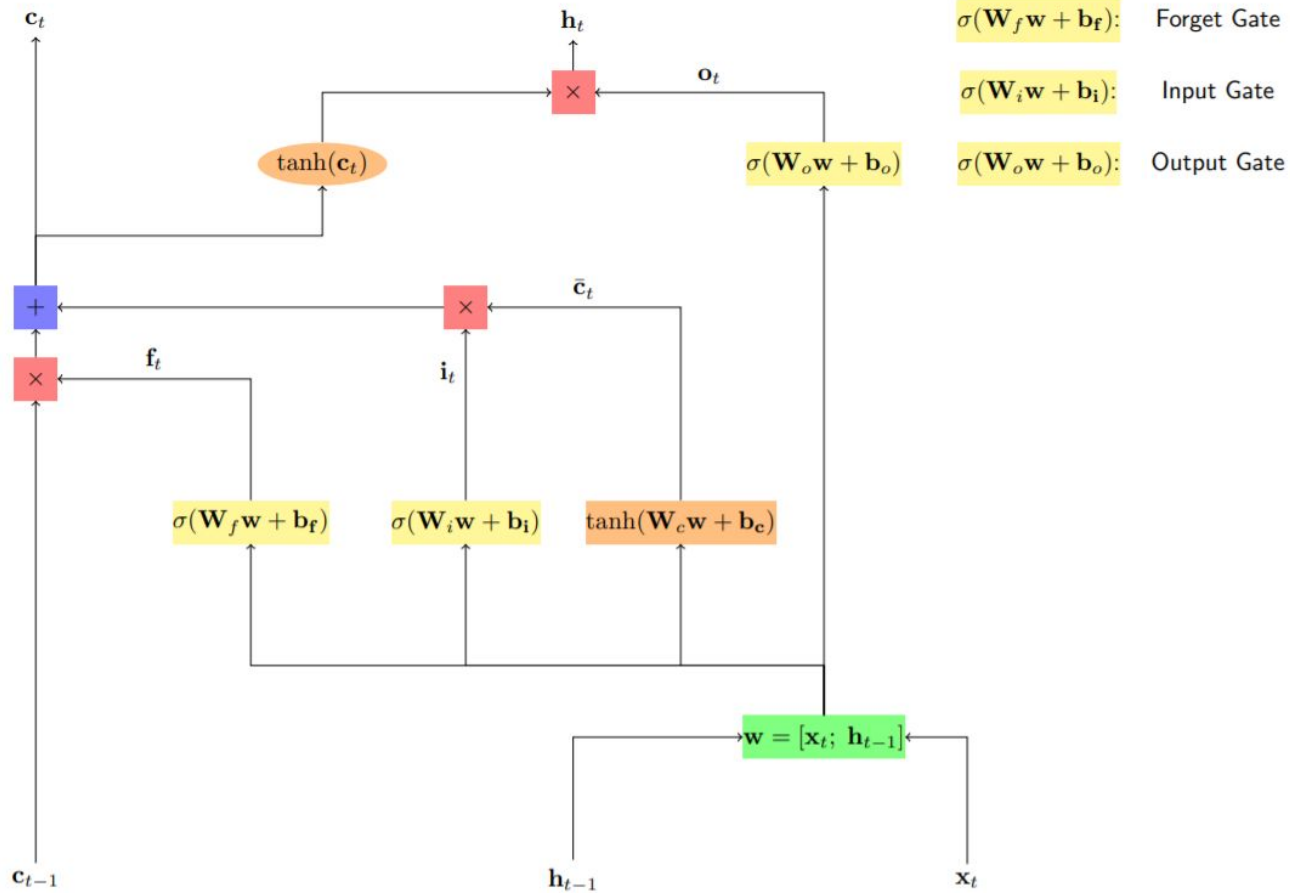


Figure: LSTM with gate equations.

The LSTM is a differentiable memory cell. We have now three gates: an input gate $\mathbf{i}^{(t)}$, a forgetting gate $\mathbf{f}^{(t)}$ and an output gate $\mathbf{o}^{(t)}$:

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

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

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

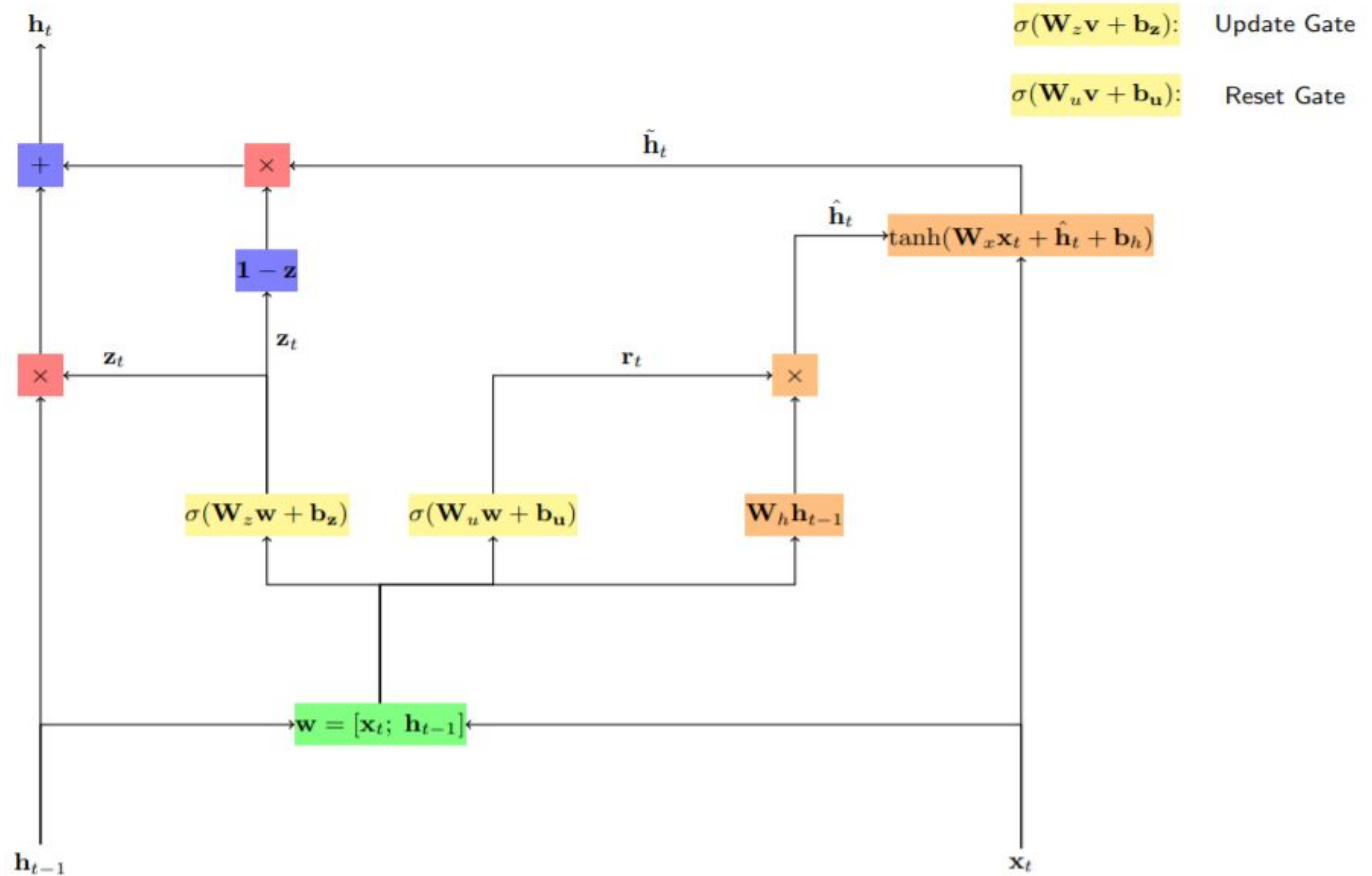
We now define a memory cell $\mathbf{c}^{(t)}$ by gating previous memory and new content

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \cdot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \cdot \tanh(\mathbf{a}^{(t)}) \quad \text{where} \quad \mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}.$$

The final hidden state is a gated activation of the memory cell:

$$\mathbf{h}_{\text{LSTM}}^{(t)} = \mathbf{o}^{(t)} \cdot \tanh \mathbf{c}^{(t)}.$$

Gated Recurrent Unit



References

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep Learning. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [2] RNNII, DLVR, Lecture by Dr. Michael Weinmann, Informatik, University of Bonn