

Recurrent Neural Networks

Deep Learning Reading Group

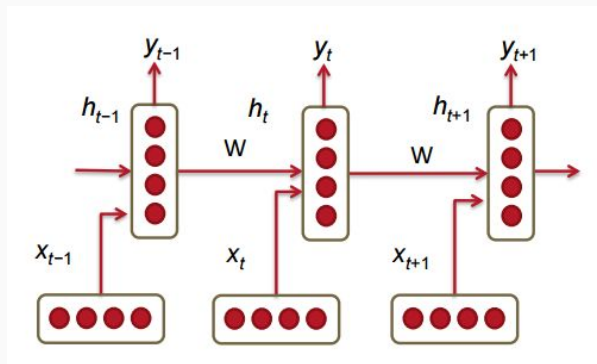
April 14, 2017

Yinong Wang



What are Recurrent Neural Networks (RNNs)?

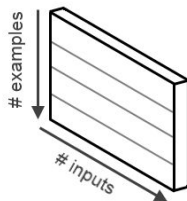
- RNNs are a type of neural network that share parameters across different positions/indices of time



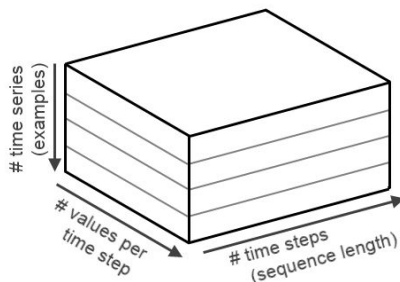
RNNs are useful for sequences

- Unlike standard NNs and CNNs, RNNs introduce a temporal element to the system
- Whereas CNNs introduce the idea of a spatial dependence, RNNs rely on the idea that the current input relies on previous inputs

Feed Forward Network Data

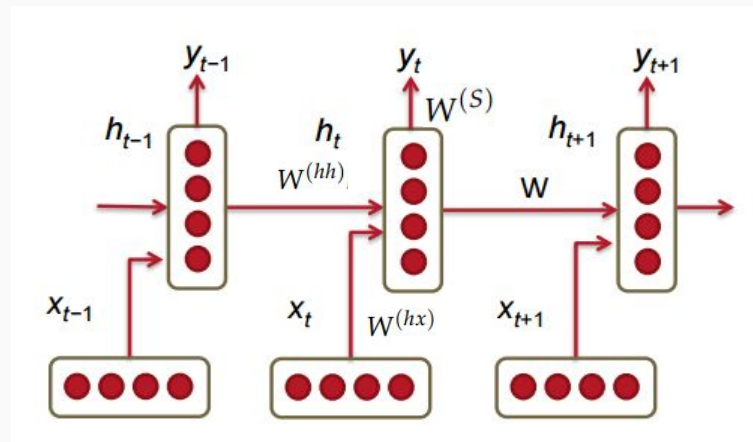


Recurrent Network Data



RNN architecture

- Current input: x_t
- Current hidden state: h_t
- Nonlinear activation function: σ , or \tanh
- Input-to-hidden weights: $W^{(hx)}$
- Hidden-to-hidden weights: $W^{(hh)}$
- Hidden-to-output weights: $W^{(S)}$
- Predicted output: \hat{y}_t



$$h_t = \sigma(W^{(hh)}h_{t-1} + W^{(hx)}x_{[t]})$$

$$\hat{y}_t = \text{softmax}(W^{(S)}h_t)$$

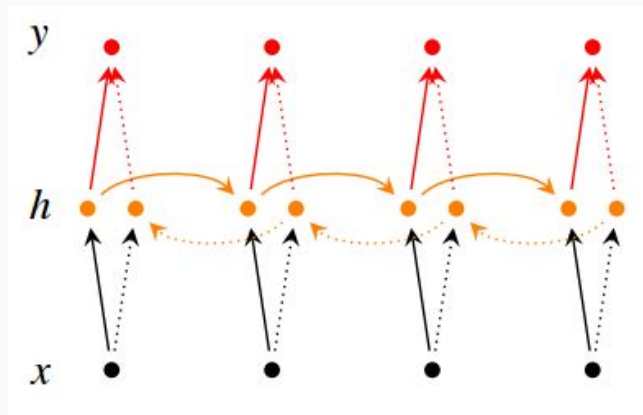
Bidirectional RNNs

- Looks at combination of previous and future inputs

$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\vec{h}_t; \overleftarrow{h}_t] + c)$$



The Vanishing (Exploding) Gradient Problem

Sentence 1

"Jane walked into the room. John walked in too. Jane said hi to John"

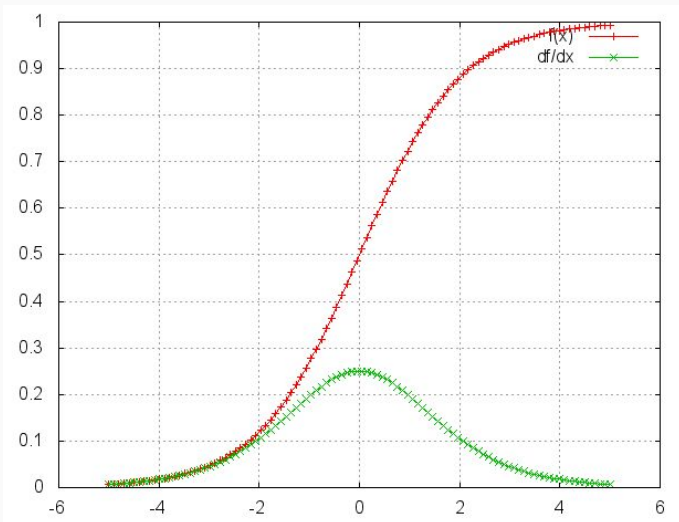
Sentence 2

"Jane walked into the room. John walked in too. It was late in the day, and everyone was walking home after a long day at work. Jane said hi to ____"

Exploring the Gradient Problem

- Saturation of neurons leads to derivatives of activation functions to be close to 0
- Gradient contribution from “far away” steps become 0
- Difficult to learn long-term dependencies
- Solutions:
 - clipping the gradient
 - ReLU activations
 - proper initialization of W
 - regularization
 - LSTMs

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \left(\prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_k}{\partial W}$$



Long Short-Term Memory units (LSTMs)

- Introduces gates: input, output, and forget
- Memory cell state

- x_t is the input to the memory cell layer at time t
- $W_i, W_f, W_c, W_o, U_i, U_f, U_c, U_o$ and V_o are weight matrices
- b_i, b_f, b_c and b_o are bias vectors

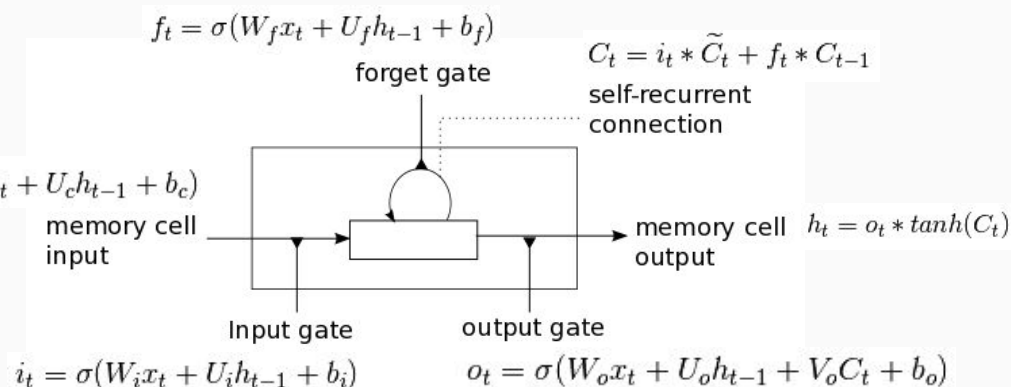


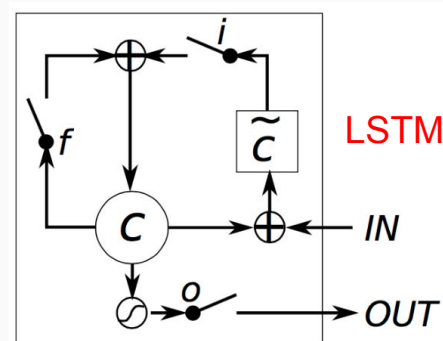
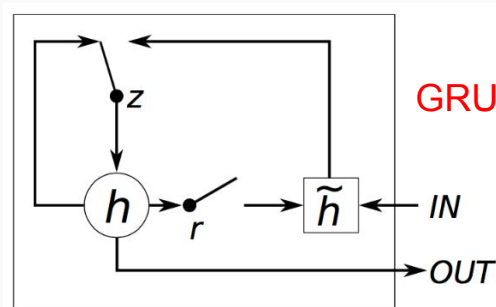
Figure 1 : Illustration of an LSTM memory cell.

Other types of RNNs

- Gated Recurrent Unit (GRU)

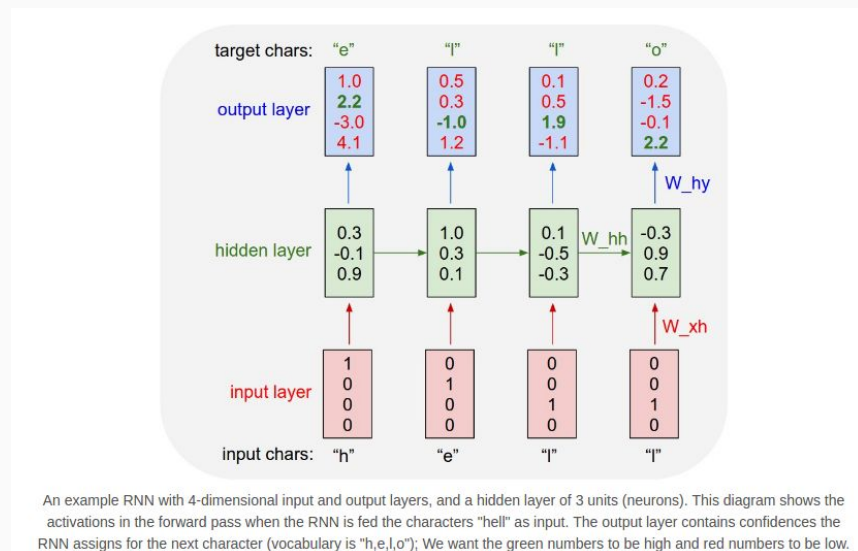
- Simplified variant of LSTM
- 2 gates, update (combination of input and output gates), and reset (similar to forget)
- No internal memory
- No additional nonlinear activation applied when computing output

$$\begin{aligned}z &= \sigma(x_t U^z + s_{t-1} W^z) \\r &= \sigma(x_t U^r + s_{t-1} W^r) \\h &= \tanh(x_t U^h + (s_{t-1} \circ r) W^h) \\s_t &= (1 - z) \circ h + z \circ s_{t-1}\end{aligned}$$



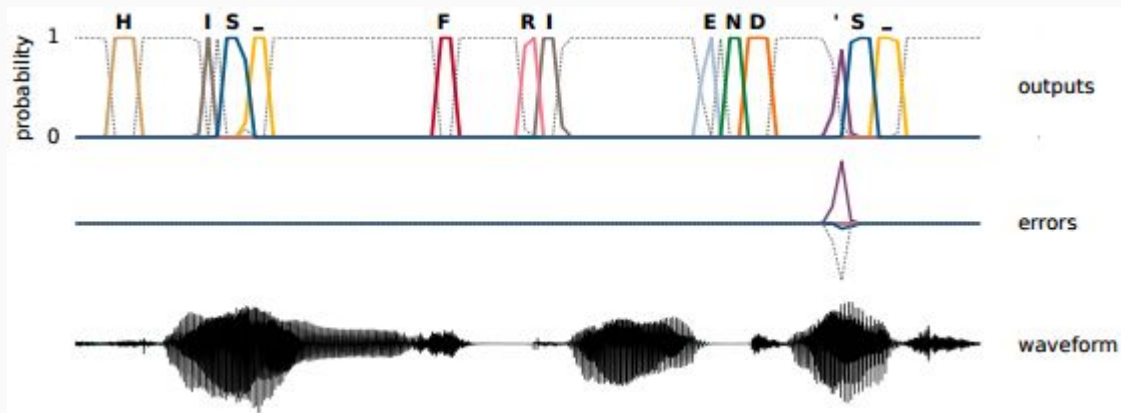
Applications of RNNs

- Character level language models



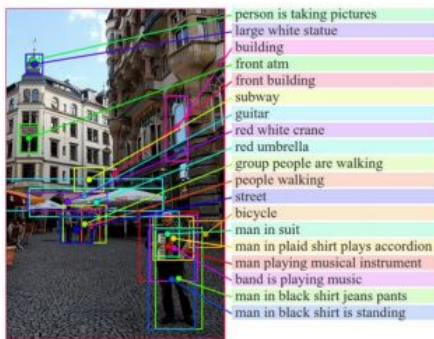
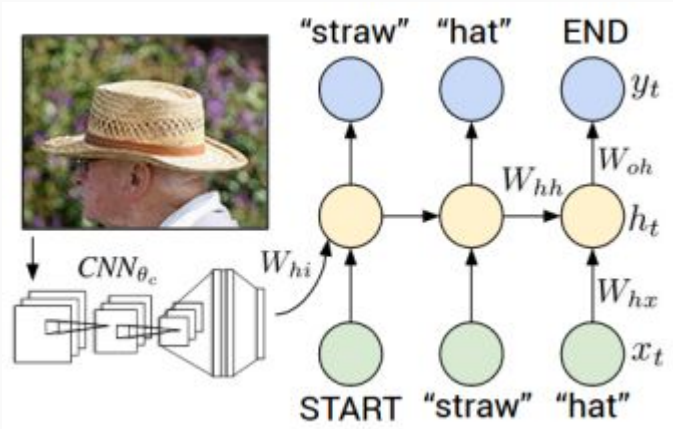
Applications of RNNs

- “Towards End-To-End Speech Recognition with Recurrent Neural Networks”
- Bidirectional LSTM



Combining RNNs and CNNs

- “Deep Visual-Semantic Alignments for Generating Image Descriptions”



construction worker in orange safety vest is working on road.



two young girls are playing with lego toy.

Trying out LSTMs

- IMDB dataset for sentiment analysis:
 - Collection of polarizing movie reviews from the Internet Movie Database (IMDB) website
 - Classification task for positive or negative movie review
 - https://github.com/fchollet/keras/blob/master/examples/imdb_lstm.py