CS60010: Deep Learning

Recurrent Neural Network

Sudeshna Sarkar

Spring 2018

5 Feb 2018

Sequence

- Sequence data: sentences, speech, stock market, signal data
 - Sequence of words in an English sentence
 - Acoustic features at successive time frames in speech recognition
 - Successive frames in video classification
 - Rainfall measurements on successive days in I'---- V-
 - Daily values of current exchange rate



Modeling Sequential Data

- Sample data sequences from a certain distribution
- Generate natural sentences to describe an image
- Activity recognition from a video sequence

 $P(y|x_1, x_2, \dots, x_T)$

 $P(x_1, x_2, ..., x_T)$

 $P(y_1, y_2, ..., y_T | I)$

• Speech Recognition

$$P(y_1, y_2, ..., y_T | x_1, x_2, ..., x_T)$$

Machine Translation

$$P(y_1, y_2, \dots, y_T | x_1, x_2, \dots, x_S)$$





Recurrent neural networks

- RNNs are very powerful, because they combine two properties:
 - Distributed hidden state that allows them to store a lot of information about the past efficiently.
 - Non-linear dynamics that allows them to update their hidden state in complicated ways.
- With enough neurons and time, RNNs can compute anything that can be computed by your computer.





Recurrent Networks offer a lot of flexibility:



Based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

Why RNNs?

- Can model sequences having variable length
- Efficient: Weights shared across time-steps

Dynamic system; Unfolded; Computation graph

Dynamical system: classical form

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

$$s^{(3)} = f(f(s^{(1)}; \theta))$$

 For a finite no. of time steps τ, the graph can be unfolded by applying the definition τ-1 times.

$$\left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(t-1)} \\ \mathbf{s}^{(t-1)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(t+1)} \\ \mathbf{s}^{(t+1)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ \mathbf{s}^{(\dots)} \\ f \end{array} \right) \xrightarrow{f} \left(\begin{array}{c} \mathbf{s}^{(\dots)} \\ \mathbf{s}^$$

The same parameters are used at all time steps.

Dynamical system driven by external signal

- Consider a dynamical system driven by external (input) signal $x^{(t)}$: $s^{(t)} = f(s^{(t-1)}, x^{(t)}; \theta)$
 - The state now contains information about the whole past input sequence. To indicate that the state is hidden rewrite using variable h for state:

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



Output prediction by RNN

- Task : To predict the future from the past
- The network typically learns to use **h**^(t) as a summary of the task-relevant aspects of the past sequence of inputs upto t
- The summary is in general lossy since it maps a sequence of arbitrary length (x^(t), x^(t-1),...,x⁽²⁾, x⁽¹⁾) to a fixed length vector h^(t)
- Depending on the training criterion, summary keeps some aspects of past sequence more precisely than other aspects

Unfolding: from circuit diagram to computational graph

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

can be written in two different ways: circuit diagram or an unfolded computational graph



The unfolded graph has a size dependent on the sequence length.

Unfolding

• We can represent the unfolded recurrence after t steps with a function $g^{(t)}$:

$$h^{(t)} = g^{(t)} (x^{(t)}, x^{(t-1)} \dots, x^{(1)})$$

- The function $g^{(t)}$ takes in whole past sequence $(x^{(t)}, x^{(t-1)} \dots, x^{(1)})$ 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. $h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta)$

Advantages of unfolding model

- Regardless of sequence length, the learned model has the same input size because it is specified in terms of transition from one state to another state rather than specified in terms of a variable length history of states
- 2. It is possible to use same function f with same parameters at every step
- These two factors make it possible to learn a single model f that operates on all time steps and all sequence lengths.
- Learning a single shared model: can apply the network to input sequences of different lengths and predict sequences of different lengths and to work with fewer training examples.

Recurrent Hidden Units



For Every Time Step!!

Recurrent Hidden Units



Unfolding Computational Graphs

 A Computational Graph is a way to formalize the structure of a set of computations such as mapping inputs and parameters to outputs and loss

 We can unfold a recursive or recurrent computation into a computational graph that has a repetitive structure



A more complex unfolded computational graph



Three design patterns of RNNs

- Output at each time step; recurrent connections between hidden units
- Output at each time step; recurrent connections only from output at one time step to hidden units at next time step
- 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



RNN1: with recurrence between hidden units

Maps input sequence **x** to output **o**

With softmax outputs Loss *L* internally computes \hat{y} = softmax(o) and compares to target *y*

• Update equation applied for each time step from t = 1 to $t = \tau$

$$egin{array}{rcl} m{a}^{(t)} &= m{b} + m{W} m{h}^{(t-1)} + m{U} m{x}^{(t)} \ m{h}^{(t)} &= ext{tanh}(m{a}^{(t)}) \ m{o}^{(t)} &= m{c} + m{V} m{h}^{(t)} \ m{y}^{(t)} &= ext{softmax}(m{o}^{(t)}) \end{array}$$



Parameters:

- bias vectors **b** and **c**
- weight matrices U (input-to-hidden),
 V (hidden-to-output) and
 W (hidden- to-hidden) connections

Loss function for a given sequence

- The total loss for a given sequence of *x* values with a sequence of *y* values is the sum of the losses over the time steps
- If L^(t) is the negative log-likelihood of y^(t) given x⁽¹⁾,..x^(t)
 then

$$L\left(\left\{x^{(1)}, x^{(2)}, \dots, x^{(t)}\right\}, \left\{y^{(1)}, y^{(2)}, \dots, y^{(t)}\right\}\right) = \sum_{t} L^{(t)}$$
$$= -\sum_{t} \log p_{model}(y^{(t)} | \{x^{(1)}, x^{(2)}, \dots, x^{(t)}\})$$

Backpropagation through time

- We can think of the recurrent net as a layered, feed-forward net with shared weights and then train the feed-forward net with weight constraints.
- We can also think of this training algorithm in the time domain:
 - The forward pass builds up a stack of the activities of all the units at each time step.
 - The backward pass peels activities off the stack to compute the error derivatives at each time step.
 - After the backward pass we add together the derivatives at all the different times for each weight.

Gradients on V, c, W and U

$$\frac{\partial L}{\partial L_t} = 1, \frac{\partial L}{\partial o_t} = \frac{\partial L}{\partial L_t} \frac{\partial L_t}{\partial o_t} = \frac{\partial L_t}{\partial o_t}$$
$$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial o_t} \frac{\partial o_t}{\partial V}$$
$$\frac{\partial L}{\partial c} = \sum_t \frac{\partial L_t}{\partial o_t} \frac{\partial o_t}{\partial c}$$
$$\frac{\partial L}{\partial W} = \sum_t \frac{\partial L_t}{\partial o_t} \frac{\partial h_t}{\partial C}$$
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h_t} \frac{\partial h_{t+1}}{\partial h_t} + \frac{\partial L}{\partial o_t} \frac{\partial o_t}{\partial h_t}$$

+ T

Feedforward Depth (d_f)



Option 2: Recurrent Depth (d_r)



Recurrent depth: Longest path between same hidden state in successive timesteps

Recurrent depth = 3