

THE UNIVERSITY OF TENNESSEE KNOXVILLE

AICIP RESEARCH

ECE599/692 - Deep Learning

Lecture 14 – Recurrent Neural Network (RNN)

Hairong Qi, Gonzalez Family Professor
 Electrical Engineering and Computer Science
 University of Tennessee, Knoxville
<http://www.eecs.utk.edu/faculty/qi>
 Email: hqi@utk.edu

AICIP RESEARCH

Outline

- Introduction
- LSTM vs. GRU
- Applications and Implementations
- References:
 - [1] Luis Serrano, A Friendly Introduction to Recurrent Neural Networks, <https://www.youtube.com/watch?v=UNmQTIOnRfg>, Aug. 2018
 - [2] Brandon Rohrer, Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM), <https://www.youtube.com/watch?v=WCUNPb-5EYI>, Jun. 2017
 - [3] Denny Britz, Recurrent Neural Networks Tutorial, <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>, Sept. 2015 (Implementation)
 - [4] Colah's blog, Understanding LSTM Networks, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>, Aug. 2015

THE UNIVERSITY OF TENNESSEE KNOXVILLE

A friendly introduction to NN [1] AICIP RESEARCH

The slide illustrates a simple neural network. On the left, a chef character is shown with a sun and a cloud. Below them is a box labeled 'NN'. On the right, two matrix equations are shown:

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

These equations correspond to the visual examples of a pizza and a burger being produced based on weather conditions (sun or rain).

A friendly introduction to RNN AICIP RESEARCH

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

THE UNIVERSITY OF TENNESSEE
4

A more complicated case AICIP RESEARCH

Monday Tuesday Wednesday Thursday Friday Saturday

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

Food

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

Same

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Next day

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Same

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Next day

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

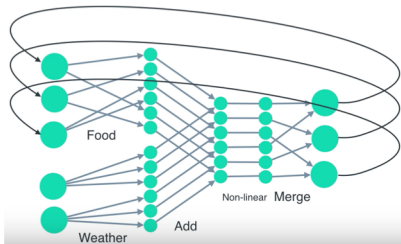
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

THE UNIVERSITY OF TENNESSEE
5

A more complicated case (cont'd) AICIP RESEARCH

THE UNIVERSITY OF TENNESSEE
6

A more complicated case (cont'd)

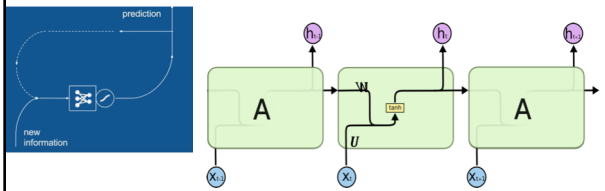


The children's book example from Brandon Rohrer [2]

Doug saw Jane.
Jane saw Spot.
Spot saw Doug.
Dictionary = {Doug, Jane, Spot, saw, .}

Potential mistakes:
Doug saw Doug.
Doug saw Jane saw Spot saw Doug...
Doug.

The standard RNN module



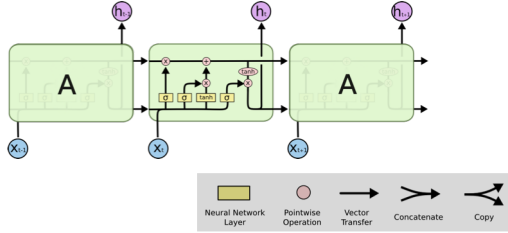
$$h_t = f(Ux_t + Wh_{t-1} + b)$$

Gradient vanishing issue: by the end of the RNN, the data from the first timestep has very little impact on the output of the RNN. An example of word prediction, "I grew up in France... I speak fluent French."

The long-short term memory (LSTM) module

AICIP RESEARCH

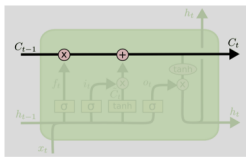
LSTMs are explicitly designed to avoid the long-term dependency problem.



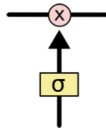
Two keys of LSTM

AICIP RESEARCH

“Cell state” which works like a conveyor belt runs straight down the entire chain, easy for information to flow along without changes.
 “Gates” which control or decide what kind of information could go or throw away from the cell state.



Cell state

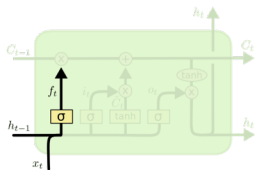
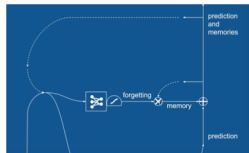


Gate

LSTM – forget gate

AICIP RESEARCH

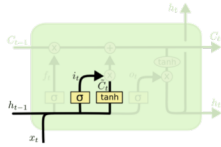
Take the example of a language model trying to predict the next word based on all the previous ones. In such a problem, the cell state might include the gender of the present subject, so that the correct pronouns can be used. When we see a new subject, we want to forget the gender of the old subject.



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTM – input gate

It decides what new information we're going to store in the cell state. It has two parts. First, a sigmoid layer called the "input gate layer" decides which values we'll update. Next, a tanh layer creates a vector of new candidate values, \tilde{C}_t , that could be added to the state. In the next step, we'll combine these two to create an update to the state.

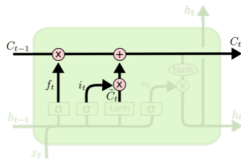


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

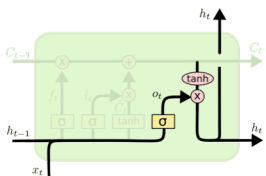
LSTM – cell state update

It actually drops the information about the old subject's gender and add the new information, as we decided in the previous steps.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTM – output gate



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

LSTM - revisit

$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$
 $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$
 $\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$
 $C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$
 $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$
 $h_t = o_t * \tanh(C_t)$

AICIP RESEARCH

THE UNIVERSITY OF TENNESSEE

16

The gated recurrent units (GRUs) module

$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$
 $r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$
 $\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$
 $h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$

Similar with LSTM but with only two gates and less parameters. The "update gate" determines how much of previous memory to be kept. The "reset gate" determines how to combine the new input with the previous memory.

AICIP RESEARCH

THE UNIVERSITY OF TENNESSEE

17

Comparison of the gating mechanism

AICIP RESEARCH

THE UNIVERSITY OF TENNESSEE

18

LSTM vs. GRU

AICIP
RESEARCH

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

THE UNIVERSITY OF
TENNESSEE

19

Application example: The talking face

AICIP
RESEARCH

Goal: Given an arbitrary audio clip and a face image, automatically generate realistic and smooth face video with accurate lip sync.

[Suwajanakorn et al., 2017]

Application: Face animation, entertainment, video bandwidth reduction, etc.

THE UNIVERSITY OF
TENNESSEE

20

The proposed framework

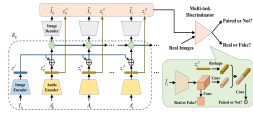
AICIP
RESEARCH

The proposed method: conditional video generation

THE UNIVERSITY OF
TENNESSEE

21

The proposed framework



$$\mathcal{L}_G = \underbrace{\mathbb{E}_{I \sim \mathcal{P}_I} [\log(1 - D(I))]}_{\text{Unconditioned}} + \underbrace{\mathbb{E}_{I \sim \mathcal{P}_I} [\log(1 - D(\tilde{I}, z_I))]}_{\text{Conditioned}} \quad (2)$$

$$\mathcal{L}_D = \underbrace{-\mathbb{E}_{I \sim \mathcal{P}_I} [\log D(I)]}_{\text{Unconditioned}} - \underbrace{\mathbb{E}_{I \sim \mathcal{P}_I} [\log(1 - D(\tilde{I}))]}_{\text{Unconditioned}} - \underbrace{\mathbb{E}_{I \sim \mathcal{P}_I} [\log(1 - D(\tilde{I}, z_I))]}_{\text{Conditioned}} \quad (3)$$

$$\mathcal{L}_{per} = \sum_{i=1}^C \underbrace{\left\| \phi_i(I) - \phi_i(\tilde{I}) \right\|}_{\text{conditioned}} \quad (4)$$

The final objective function is

$$\min_{G, E_A, E_I, R_E} \max_D \mathcal{F}(G, E_A, E_I, R_E, D), \quad (5)$$

where $\mathcal{F}(G, E_A, E_I, R_E, D) = L_{rec} + \gamma L_{per} + \lambda L_G + \lambda L_D$.