Recurrent Neural Networks I

CS 4391 Introduction Computer Vision Professor Yu Xiang The University of Texas at Dallas

Some slides of this lecture are courtesy Stanford CS231n

NIN

Single Images

• Convolutional neural networks



High-level information

- Depth
- Object classes
- Object poses
- Etc.

Sequential Data

- Data depends on time
 - Video



• Sentence

UT Dallas is a rising public research university in the heart of DFW.

Sequential Data Labeling

• Video frame labeling



Frames of a Video

https://bleedai.com/human-activity-recognition-using-tensorflow-cnn-lstm/

Yu Xiang

Sequential Data Labeling

Part-of-speech tagging (grammatical tagging)



Tag	Meaning	English Examples
ADJ	adjective	new, good, high, special, big, local
ADP	adposition	on, of, at, with, by, into, under
ADV	adverb	really, already, still, early, now
CONJ	conjunction	and, or, but, if, while, although
DET	determiner, article	the, a, some, most, every, no, which
NOUN	noun	year, home, costs, time, Africa
NUM	numeral	twenty-four, fourth, 1991, 14:24
PRT	particle	at, on, out, over per, that, up, with
PRON	pronoun	he, their, her, its, my, I, us
VERB	verb	is, say, told, given, playing, would
	punctuation marks	.,;!
х	other	ersatz, esprit, dunno, gr8, univeristy

Sequential Data Labeling



Recurrent Neural Networks



Hidden State Update



Using the Hidden State



Recurrent Neural Networks



Vanilla RNN



RNN Computation Graph



The same set of weights for different time steps $\ f_W \ f_{W'}$



Backpropagation through Time



Yu Xiang

Truncated Backpropagation through Time

Run forward and backward through chunks of the sequence instead of whole sequence

Truncated Backpropagation through Time

Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

Yu Xiang

Truncated Backpropagation through Time

$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W}$$
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \dots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} (\prod_{t=2}^{T} \frac{\partial h_t}{\partial h_{t-1}}) \frac{\partial h_1}{\partial W}$$

$$rac{\partial L_T}{\partial W} = rac{\partial L_T}{\partial h_T} (\prod_{t=2}^T rac{\partial h_t}{\partial h_{t-1}}) rac{\partial h_1}{\partial W}$$

https://en.wikipedia.org/wiki/Matrix_norm

 Vanishing gradients $\|\frac{\partial h_t}{\partial h_{t-1}}\|_2 < 1$

 Exploding gradients

• Exploding gradients

$$\|\frac{\partial h_t}{\partial h_{t-1}}\|_2 > 1$$

01

- Gradient clipping
- grad_norm = np.sum(grad * grad)
 if grad_norm > threshold:
 grad *= (threshold / grad_norm)
- Vanishing gradients

$$\frac{\partial h_t}{\partial h_{t-1}} \|_2 < 1$$

• Change RNN architecture

Summary

- RNNs can be used for sequential data to capture dependencies in time
- LSTMs and GRUs are better then vanilla RNNs
- It is difficult to capture long-term dependencies in RNNs
- Use transformers (in future lectures)

Further Reading

- Stanford CS231n, lecture 10, Recurrent Neural Networks <u>http://cs231n.stanford.edu/</u>
- Long Short Term Memory <u>https://www.researchgate.net/publication/13853244 Long Short-</u> <u>term Memory</u>
- Gated Recurrent Units <u>https://arxiv.org/pdf/1412.3555.pdf</u>