CS388: Natural Language Processing Lecture 8: RNNs



Greg Durrett





Project 1 due Thursday at 5pm

Administrivia



- heuristics (e.g., Glorot initializer)
- Dropout is an effective regularizer

Think about your optimizer: Adam or tuned SGD work well



Recall: Training Tips

Parameter initialization is critical to get good gradients, some useful

Generative Parsing (Training Set)

Generative Parsing (Development Set) (f)



Recall: Word Vectors

	theore	sident said tha	it the downturn was o
	president	the of	
	president	the said	president
	governor	the of	governor
	governor	the appointed	
	said	sources •	said)
	said	president that	reportea
	reported	sources	

[Finch and Chater 92, Shuetze 93, many others]





Recall: Continuous Bag-of-Words

Predict word from context





Matrix factorization approaches useful for learning





Using Word Embeddings

- no pretraining
 - Often works pretty well
- Approach 2: pretrain using GloVe, keep fixed
 - Faster because no need to update these parameters
 - Need to make sure GloVe vocabulary contains all the words you need
- Approach 3: initialize using GloVe, fine-tune
 - Not as commonly used anymore

Approach 1: learn embeddings directly from data in your neural model,





- What if we want embedding representations for whole sentences?
- Skip-thought vectors (Kiros et al., 2015), similar to skip-gram generalized to a sentence level (more later)
- Is there a way we can compose vectors to make sentence representations? Summing? RNNs?

Compositional Semantics





Recurrent neural networks

Vanishing gradient problem

LSTMs / GRUs

Applications / visualizations

This Lecture



RNN Motivation



Feedforward NNs can't handle variable length input: each position in the feature vector has fixed semantics



- These don't look related (great is in two different orthogonal subspaces)
- Instead, we need to:
- 1) Process each word in a uniform way
- 2) ... while still exploiting the context that that token occurs in







RNN Abstraction



hidden state and produces output y (all vector-valued)



Cell that takes some input x, has some hidden state h, and updates that



Transducer: make some prediction for each element in a sequence



Acceptor/encoder: encode a sequence into a fixed-sized vector and use that for some purpose



- output **y** = score for each tag, then softmax

- predict sentiment (matmul + softmax)
- translate
- paraphrase/compress







Long history! (invented in the late 1980s)

$$\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_t)$$

Updates hidden state based on input and current hidden state

$$\mathbf{y}_t = \tanh(U\mathbf{h}_t + \mathbf{b}_y)$$

Computes output from hidden state









Training Elman Networks



- Backpropagation through time": build the network as one big computation graph, some parameters are shared
- RNN potentially needs to learn how to "remember" information for a long time!
- it was my favorite movie of 2016, though it wasn't without problems -> +
- "Correct" parameter update is to do a better job of remembering the sentiment of *favorite*

predict sentiment





Vanishing Gradient



LSTMs/GRUs



- Designed to fix "vanishing gradient" problem using gates $\mathbf{h}_t = \mathbf{h}_{t-1} \odot \mathbf{f} + \operatorname{func}(\mathbf{x}_t)$ gated
- Vector-valued "forget gate" f computed based on input and previous hidden state

$$\mathbf{f} = \sigma(W^{xf}\mathbf{x}_t + W^{hf}\mathbf{h}_{t-1})$$

- Sigmoid: elements of f are in [0, 1]
- If f = 1, we simply sum up a function of all inputs — gradient doesn't vanish!

Gated Connections

$\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_h)$ Elman







LSTMS

"Cell" c in addition to hidden state h $\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f} + \operatorname{func}(\mathbf{x}_t, \mathbf{h}_{t-1})$

Vector-valued forget gate f depends on the h hidden state

$$\mathbf{f} = \sigma(W^{xf}\mathbf{x}_t + W^{hf}\mathbf{h}_{t-1})$$

in addition to gates from previous timesteps

Basic communication flow: x -> c -> h, each step of this process is gated





LSTMs

 $\mathbf{c_{j}=}\mathbf{c_{j-1}}\odot\mathbf{f}+\mathbf{g}\odot\mathbf{i}$ $\mathbf{f} = \sigma(\mathbf{x_j}\mathbf{W^{xf}} + \mathbf{h_{j-1}}\mathbf{W^{hf}})$ $\mathbf{g} = \operatorname{tanh}(\mathbf{x_j} \mathbf{W^{xg}} + \mathbf{h_{j-1}} \mathbf{W^{hg}})$ $\mathbf{i} = \sigma(\mathbf{x_j} \mathbf{W^{xi}} + \mathbf{h_{j-1}} \mathbf{W^{hi}})$ $\mathbf{h_i} = \operatorname{tanh}(\mathbf{c_i}) \odot \mathbf{o}$ $\mathbf{o} = \sigma(\mathbf{x_i} \mathbf{W^{xo}} + \mathbf{h_{i-1}} \mathbf{W^{ho}})$

Goldberg lecture notes









LSTMs

 $\mathbf{c_{j}=}\mathbf{c_{j-1}}\odot\mathbf{f}+\mathbf{g}\odot\mathbf{i}$ $\mathbf{f} = \sigma(\mathbf{x_j}\mathbf{W^{xf}} + \mathbf{h_{j-1}}\mathbf{W^{hf}})$ $\mathbf{g} = \operatorname{tanh}(\mathbf{x_j}\mathbf{W^{xg}} + \mathbf{h_{j-1}}\mathbf{W^{hg}})$ $\mathbf{i} = \sigma(\mathbf{x_j} \mathbf{W^{xi}} + \mathbf{h_{j-1}} \mathbf{W^{hi}})$

 $\mathbf{h_i} = \operatorname{tanh}(\mathbf{c_i}) \odot \mathbf{o}$

 $\mathbf{o} = \sigma(\mathbf{x_j} \mathbf{W^{xo}} + \mathbf{h_{j-1}} \mathbf{W^{ho}})$

Can we ignore c in our current computation?

Can we output something without changing c?









LSTMs

- Ignoring recurrent state entirely:
 - Lets us get feedforward layer over token
- Ignoring input:
 - Lets us discard stopwords
- Summing inputs:
 - Lets us compute a bag-of-words representation

Goldberg lecture notes









usually initialize forget gate = 1 to remember everything to start

LSTMs

Gradient still diminishes, but in a controlled way and generally by less —



Understanding LSTM Parameters





Initialize hidden layer randomly

Need to learn how the gates work: what do we forget/remember?

g uses an arbitrary nonlinearity, this is the "layer" of the cell













LSTM: more complex and slower, may work a bit better

GRUS



GRU: faster, a bit simpler

Two gates: **z** (forget, mixes **s** and h) and r (mixes h and x)

What do RNNs produce?





- Encoding of the sentence can pass this a decoder or make a classification decision about the sentence
- Encoding of each word can pass this to another layer to make a prediction (can also pool these to get a different sentence encoding)
- RNN can be viewed as a transformation of a sequence of vectors into a sequence of context-dependent vectors

the movie was great



Multilayer Bidirectional RNN



Sentence classification
based on concatenation
of both final outputs



Token classification based on concatenation of both directions' token representations



Training RNNs





the movie was great

- Loss = negative log likelihood of probability of gold label (or use SVM) or other loss)
- Backpropagate through entire network
- Example: sentiment analysis

Training RNNs





the movie was great

- Loss = negative log likelihood of probability of gold predictions, summed over the tags
- Loss terms filter back through network
- Example: language modeling (predict next word given context)

Applications



What can LSTMs model?

- Sentiment
 - Encode one sentence, predict
- Language models
 - Move left-to-right, per-token prediction
- Translation
 - Encode sentence + then decode, use token predictions for attention weights (later in the course)



- Counter: know when to generate \n

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae-pressed forward into boats and into the ice-covered water and did not, surrender.

Train character LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code

Visualize activations of specific cells (components of c) to understand them









- Visualize activations of specific cells to see what they track
- Binary switch: tells us if we're in a quote or not



Train character LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code





- Train character LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- Visualize activations of specific cells to see what they track
- Stack: activation based on indentation







- Train character LSTM language model (predict next character based on history) over two datasets: War and Peace and Linux kernel source code
- Visualize activations of specific cells to see what they track
- Uninterpretable: probably doing double-duty, or only makes sense in the context of another activation











What can LSTMs model?

- Sentiment
 - Encode one sentence, predict
- Language models
 - Move left-to-right, per-token prediction
- Translation
 - Encode sentence + then decode, use token predictions for attention weights (next lecture)
- Textual entailment
 - Encode two sentences, predict



Natural Language Inference

Premise

A boy plays in the snow

A man inspects the uniform of a figure

An older and younger man smiling

- 2006 (Dagan, Glickman, Magnini)
- knowledge, temporal reasoning, etc.)

Hypothesis

entails A boy is outside

The man is sleeping contradicts Two men are smiling and neutral laughing at cats playing

Long history of this task: "Recognizing Textual Entailment" challenge in

Early datasets: small (hundreds of pairs), very ambitious (lots of world





- contradictory statements
- >500,000 sentence pairs
- Encode each sentence and process 100D LSTM: 78% accuracy 300D LSTM: 80% accuracy Bowman et al., 2016) 300D BiLSTM: 83% accuracy (Liu et al., 2016) Later: better models for this

SNLI Dataset

Show people captions for (unseen) images and solicit entailed / neural /



Bowman et al. (2015)





- RNNs can transduce inputs (produce one output for each input) or compress the whole input into a vector
- Useful for a range of tasks with sequential input: sentiment analysis, language modeling, natural language inference, machine translation
- Next time: CNNs and neural CRFs