

# EECS 230 Deep Learning Lecture 11: RNN and LSTM

Credit to Daniel Jurafsky for figures of network architectures

# Recap: Language models

#### Many NLP tasks require **natural language output**:

- -Machine translation
- -Speech recognition
- -Natural language generation
- -Spell-checking

#### Language models define **probability distributions over** (natural language) **strings or sentences**.

- → We can use a language model to generate strings
- → We can use a language model to score/rank candidate strings so that we can choose the best (i.e. most likely) one:

if PLM(A) > PLM(B), return A, not B



# Recap: N-gram Language model

N-gram models *assume* each word (event) depends only on the previous n–1 words (events): Unigram model:  $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^{N} P(w^{(i)})$ Bigram model:  $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^{N} P(w^{(i)} | w^{(i-1)})$ Trigram model:  $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^{N} P(w^{(i)} | w^{(i-1)}, w^{(i-2)})$ 

Independence assumptions where the n-th event in a sequence depends only on the last n-1 events are called Markov assumptions (of order n-1).



#### Recap: Learning a language model

We need (a large amount of) text as training data to estimate the parameters of a language model.

The most basic parameter estimation technique: relative frequency estimation (frequency = counts)  $P(w^{(i)} = 'the' | w^{(i-1)} = 'on') = C('on the') / C('on')$ Also called Maximum Likelihood Estimation (MLE)



#### **Recap: Word Embedding**

expect =



been

had have

need

help



#### Recap: Word2Vec Overview

**D**Example windows and process for computing  $P(w_{t+j} | w_t)$ 





#### Recap: Word2Vec Objective Function

For each position t = 1, ..., T, predict context words within a window of fixed size m, given center word  $w_i$ . Data likelihood:

Likelihood = 
$$L(\theta) = \prod_{t=1}^{T} \prod_{\substack{-m \le j \le m \\ j \ne 0}} P(w_{t+j} \mid w_t; \theta)$$
  
 $\theta$  is all variables  
to be optimized  
sometimes called a *cost* or *loss* function

The objective function  $J(\theta)$  is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T}\log L(\theta) = -\frac{1}{T}\sum_{t=1}^{T}\sum_{\substack{-m \le j \le m \\ j \ne 0}}\log P(w_{t+j} \mid w_t; \theta)$$

Minimizing objective function ⇔ Maximizing predictive accuracy



# Our first neural net for NLP: A neural n-gram model

Given a fixed-size vocabulary V, an *n*-gram model predicts the probability of the *n*-th word following the preceding *n*–1 words:

$$P(w^{(i)} | w^{(i-1)}, w^{(i-2)}, \dots, w^{i-(n-1)})$$

How can we model this with a neural net?

- Input layer: concatenate n-1 word vectors
- Output layer: a softmax over IVI units
- Feedforward network



# Outline

Recurrent Neural Network

**RNNs** as Language Models

□RNNs for other NLP tasks

□Staked and Bi-directional RNN

Encoder-decoder model



#### **Recurrent Neural Network**

Temporal nature in language processing

**Q**RNN deals with sequential input data stream like language.



A simple RNN



# A Simple Recurrent Neural Network

**Q**RNN illustrated as a feed-forward network





# A Simple Recurrent Neural Network

#### **RNN** unrolled in time



#### How to optimize Recurrent Neural Network?



□ Backpropagation through time

$$\frac{\partial L_3}{\partial \mathbf{W}} = \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{W}} + \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}} + \frac{\partial L_3}{\partial \mathbf{h}_3} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}}$$

$$\frac{\partial L}{\partial \mathbf{W}} = -\frac{1}{n} \sum_{t=1}^{n} \sum_{k=1}^{t} \frac{\partial L_t}{\partial \mathbf{h}_t} \left( \prod_{j=k+1}^{t} \frac{\partial \mathbf{h}_j}{\partial \mathbf{h}_{j-1}} \right) \frac{\partial \mathbf{h}_k}{\partial \mathbf{W}}$$



# Truncated backpropagation through time

□ Backpropagation is very expensive for long sequences



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



#### RNNs as Language Models

Language models predict the next word in a sequence given some preceding context.

P(fish|Thanks for all the)

RNN Language Model



# Training an RNN Language Model

#### Maximum likelihood estimation





# Generation with RNN Language Model

Autoregressive (casual) generation

□Unlike teacher forcing during training



#### RNNs for other NLP tasks

□RNN for sequence classification

Commonly called text classification, like sentiment analysis or spam detection



#### Stacked RNN

Outputs from one RNN as an input sequence to another one



### **Bidirectional RNN**

Two independent RNN, one forward, another backward



# Bidirectional RNN for sequence classification

Combine hidden states from first & last word



# Vanishing/exploding gradients

Consider the gradient of  $L_t$  at step t, with respect to the hidden state  $\mathbf{h}_k$  at some previous step k (k<t ):

$$\frac{\partial L_t}{\partial \mathbf{h}_k} = \frac{\partial L_t}{\partial \mathbf{h}_t} \left( \prod_{t \ge j > k} \frac{\partial \mathbf{h}_j}{\partial \mathbf{h}_{j-1}} \right)$$

□ Recurrent multiplication

Gradients too small (vanishing gradient) or too large (exploding gradient)



# **Exploding gradients**

**What is the problem?** 

□We take a very large step in SGD

□Solution: Gradient clipping





# Vanishing gradients

- **W**hat is the problem?
- □ Parameters barely get updated (no learning)

□Solution:

- □LSTMs: Long short-term memory networks
- GRUs: Gated recurrent units



# Local vs distant information

Hidden states tend to contain local information
 But distant information is critical

"The flights the airline was canceling were full" Should predict "were" given distant information (flights)



# Long Short-term Memory (LSTM)

A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem
 Basic idea: turning multiplication into addition
 Use "gates" to control how much information to add/erase
 At each timestep, there is a hidden state h<sub>t</sub> (local information) and also a cell state C<sub>t</sub> (distant information)



# Long Short-term Memory (LSTM)

- Gate: feedforward layer, followed by a sigmoid activation function, followed by a pointwise multiplication with the layer being gated
- Generation For example, output gate (What to output for hidden state)

$$\mathbf{o}_{t} = \boldsymbol{\sigma}(\mathbf{U}_{o}\mathbf{h}_{t-1} + \mathbf{W}_{o}\mathbf{x}_{t})$$
  
$$\mathbf{h}_{t} = \mathbf{o}_{t} \odot \tanh(\mathbf{c}_{t})$$

#### □Other gates

- □Forget gate
- Add gate
- □Input gate



#### Long Short-term Memory (LSTM)





### Summary: Common RNN NLP architectures



# Encoder Decoder Architecture

□Arbitrary length output given an input sequence

- **D**A.K.A. sequence-to-sequence network
- Context vector conveys the essence of the input to the decoder





## Encoder Decoder Architecture

#### □Training an encoder-decoder for machine translation





# Problem of Encoder-decoder architecture

Context vector encodes EVERYTHING about input sequenceContext vector acts as a bottleneck





#### **Attention Mechanism**

Each output in decoder accesses all the hidden states from the encoder, not just the last state

Each output attends to all input



# What's next?

- □Self-attention
- □Transformer block
- □Transformer architecture
- Large Language Model
- □Vision Transformers

