# EE-608: Deep Learning For NLP: Language Modelling, Recurrent Neural Networks, LSTMs

James Henderson



DLNLP, Lecture 2

# **Outline**

Language Modelling

Recurrent Neural Networks

Deep Backprop with RNNs (LSTMs)

Bigger is Better with RNNs

# Background on Neural Networks and Derivatives

Everyone should make sure they know the material covered in cs224n-2023-lecture03-neuralnets.pdf at

http://web.stanford.edu/class/cs224n/

# **Outline**

# Language Modelling

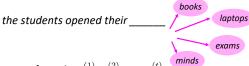
Recurrent Neural Networks

Deep Backprop with RNNs (LSTMs

Bigger is Better with RNNs

# 2. Language Modeling

Language Modeling is the task of predicting what word comes next



• More formally: given a sequence of words  $x^{(1)}, x^{(2)}, \dots, x^{(t)}$ , compute the probability distribution of the next word  $x^{(t+1)}$ :

$$P(x^{(t+1)}|x^{(t)},...,x^{(1)})$$

where  $oldsymbol{x}^{(t+1)}$  can be any word in the vocabulary  $V = \{oldsymbol{w}_1, ..., oldsymbol{w}_{|V|}\}$ 

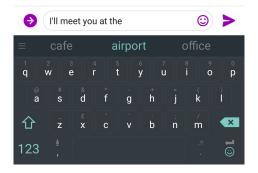
A system that does this is called a Language Model

### **Language Modeling**

- You can also think of a Language Model as a system that assigns a probability to a piece of text
- For example, if we have some text  $x^{(1)},\dots,x^{(T)}$  , then the probability of this text (according to the Language Model) is:

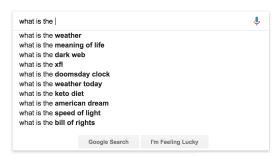
$$\begin{split} P(\boldsymbol{x}^{(1)},\dots,\boldsymbol{x}^{(T)}) &= P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)}|\ \boldsymbol{x}^{(1)}) \times \dots \times P(\boldsymbol{x}^{(T)}|\ \boldsymbol{x}^{(T-1)},\dots,\boldsymbol{x}^{(1)}) \\ &= \prod_{t=1}^T P(\boldsymbol{x}^{(t)}|\ \boldsymbol{x}^{(t-1)},\dots,\boldsymbol{x}^{(1)}) \end{split}$$

# You use Language Models every day!



# You use Language Models every day!





### n-gram Language Models

the students opened their \_\_\_\_\_

- Question: How to learn a Language Model?
- Answer (pre- Deep Learning): learn an n-gram Language Model!
- **Definition:** An *n*-gram is a chunk of *n* consecutive words.
  - · unigrams: "the", "students", "opened", "their"
  - · bigrams: "the students", "students opened", "opened their"
  - · trigrams: "the students opened", "students opened their"
  - · four-grams: "the students opened their"
- Idea: Collect statistics about how frequent different n-grams are and use these to predict next word.

### n-gram Language Models

• First we make a Markov assumption:  $x^{(t+1)}$  depends only on the preceding n-1 words

$$P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)},\dots,\boldsymbol{x}^{(1)}) = P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)},\dots,\boldsymbol{x}^{(t-n+2)})$$
 (assumption) prob of a n-gram 
$$P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)},\dots,\boldsymbol{x}^{(t-n+2)})$$
 (definition of conditional prob)

- Question: How do we get these *n*-gram and (*n*-1)-gram probabilities?
- Answer: By counting them in some large corpus of text!

$$pprox rac{\mathrm{count}(oldsymbol{x}^{(t+1)},oldsymbol{x}^{(t)},\ldots,oldsymbol{x}^{(t-n+2)})}{\mathrm{count}(oldsymbol{x}^{(t)},\ldots,oldsymbol{x}^{(t-n+2)})}$$
 (statistical approximation)

### n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

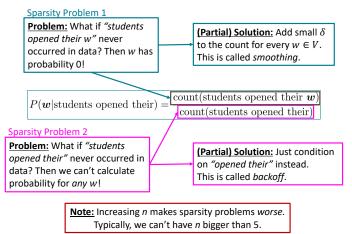
as the proctor started the clock, the students opened their discard condition on this 
$$P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their }\boldsymbol{w})}{\text{count}(\text{students opened their})}$$

For example, suppose that in the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

# **Sparsity Problems with n-gram Language Models**



# **Storage Problems with n-gram Language Models**

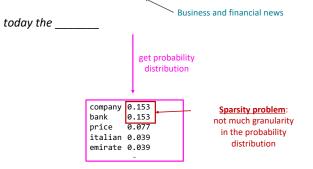
**Storage**: Need to store count for all *n*-grams you saw in the corpus.

 $P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w})}{\text{count}(\text{students opened their})}$ 

Increasing *n* or increasing corpus increases model size!

### n-gram Language Models in practice

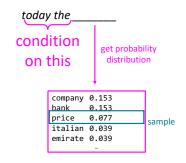
You can build a simple trigram Language Model over a
 1.7 million word corpus (Reuters) in a few seconds on your laptop\*



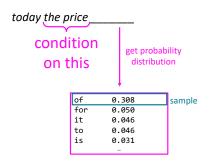
Otherwise, seems reasonable!

<sup>\*</sup> Try for yourself: https://nlpforhackers.io/language-models/

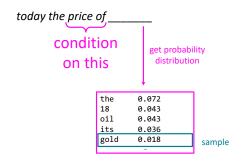
You can also use a Language Model to generate text



You can also use a Language Model to generate text



You can also use a Language Model to generate text



You can also use a Language Model to generate text

today the price of gold per ton, while production of shoe lasts and shoe industry, the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks, sept 30 end primary 76 cts a share.

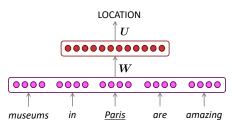
Surprisingly grammatical!

...but **incoherent.** We need to consider more than three words at a time if we want to model language well.

But increasing *n* worsens sparsity problem, and increases model size...

# How to build a neural language model?

- Recall the Language Modeling task:
  - Input: sequence of words  $x^{(1)}, x^{(2)}, \dots, x^{(t)}$
  - Output: prob. dist. of the next word  $P({m x}^{(t+1)}|\ {m x}^{(t)},\dots,{m x}^{(1)})$
- How about a window-based neural model?
  - We saw this applied to Named Entity Recognition in Lecture 2:



# A fixed-window neural Language Model



# A fixed-window neural Language Model

output distribution

$$\hat{\boldsymbol{y}} = \operatorname{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

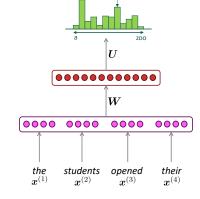
hidden laver

$$h = f(We + b_1)$$

concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

words / one-hot vectors  $oldsymbol{x}^{(1)}, oldsymbol{x}^{(2)}, oldsymbol{x}^{(3)}, oldsymbol{x}^{(4)}$ 



books

laptops

# A fixed-window neural Language Model

Approximately: Y. Bengio, et al. (2000/2003): A Neural Probabilistic Language Model

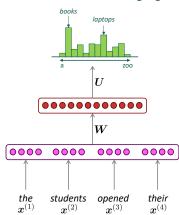
#### **Improvements** over *n*-gram LM:

- No sparsity problem
- Don't need to store all observed n-grams

#### Remaining problems:

- · Fixed window is too small
- Enlarging window enlarges W
- Window can never be large enough!
- x<sup>(1)</sup> and x<sup>(2)</sup> are multiplied by completely different weights in W.
   No symmetry in how the inputs are processed.

We need a neural architecture that can process any length input



# **Outline**

Language Modelling

**Recurrent Neural Networks** 

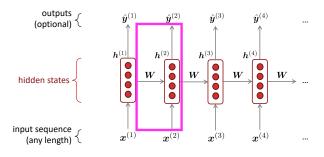
Deep Backprop with RNNs (LSTMs)

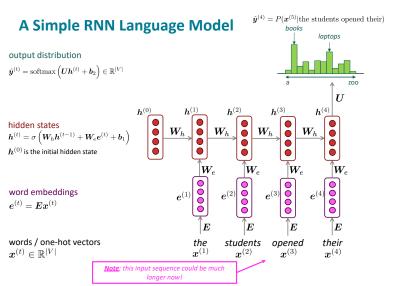
Bigger is Better with RNNs

#### 3. Recurrent Neural Networks (RNN)

A family of neural architectures

Core idea: Apply the same weights *W repeatedly* 





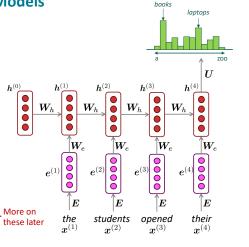
# **RNN Language Models**

#### RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

#### **RNN Disadvantages:**

- Recurrent computation is slow
- In practice, difficult to access information from many steps back



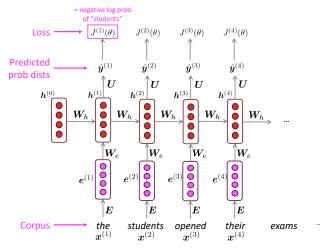
 $\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$ 

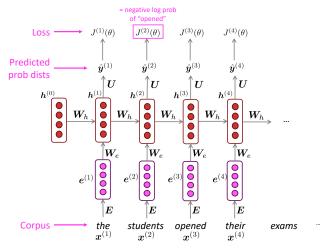
- Get a big corpus of text which is a sequence of words  $x^{(1)}, \dots, x^{(T)}$
- Feed into RNN-LM; compute output distribution  $\hat{y}^{(t)}$  for every step t.
  - · i.e., predict probability dist of every word, given words so far
- Loss function on step t is cross-entropy between predicted probability distribution  $\hat{y}^{(t)}$ , and the true next word  $y^{(t)}$  (one-hot for  $x^{(t+1)}$ ):

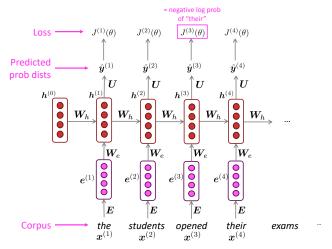
$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

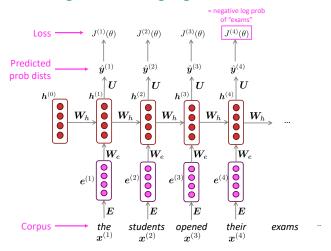
• Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

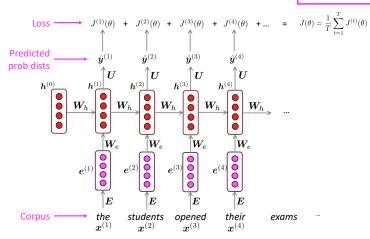








"Teacher forcing"

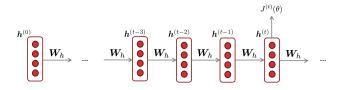


• However: Computing loss and gradients across entire corpus  $x^{(1)}, \dots, x^{(T)}$  at once is too expensive (memory-wise)!

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta)$$

- In practice, consider  $m{x}^{(1)},\dots,m{x}^{(T)}$  as a sentence (or a document)
- Recall: Stochastic Gradient Descent allows us to compute loss and gradients for small chunk of data, and update.
- Compute loss  $J(\theta)$  for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat on a new batch of sentences.

# **Backpropagation for RNNs**



**Question:** What's the derivative of  $J^{(t)}(\theta)$  w.r.t. the repeated weight matrix  $W_h$  ?

$$\begin{array}{ll} \underline{\text{Answer:}} & \frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} = \sum_{i=1}^t \left. \frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} \right|_{(i)} \end{array}$$

"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

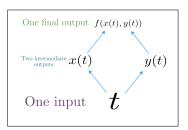
Why?

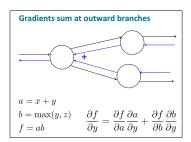
#### **Multivariable Chain Rule**

- Given a multivariable function f(x,y), and two single variable functions x(t) and y(t), here's what the multivariable chain rule says:

$$rac{d}{dt} f(x(t), extbf{ extit{y}}(t)) = rac{\partial f}{\partial x} rac{dx}{dt} + rac{\partial f}{\partial extbf{ extit{y}}} rac{dy}{dt}$$

Derivative of composition function

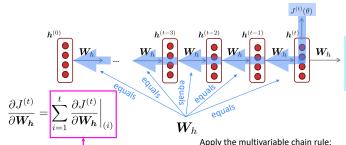




#### Source:

https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-rule-simple-version

# Training the parameters of RNNs: Backpropagation for RNNs



In practice, often "truncated" after ~20 timesteps for training efficiency reasons

Question: How do we calculate this?

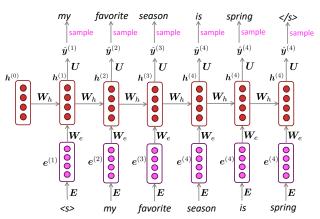
Answer: Backpropagate over timesteps i = t, ..., 0, summing gradients as you go. This algorithm is called "backpropagation through time" [Werbos, P.G., 1988, Neural Networks 1, and others]

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}} \Big|_{(i)} \frac{\partial \boldsymbol{W}_{h}|_{(i)}}{\partial \boldsymbol{W}_{h}}$$
$$= \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}} \Big|_{(i)}$$

Slide from Christopher Manning

## Generating with an RNN Language Model ("Generating roll outs")

Just like an n-gram Language Model, you can use a RNN Language Model to generate text by repeated sampling. Sampled output becomes next step's input.



## **Generating text with an RNN Language Model**

#### Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on Obama speeches:



The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

Source: https://medium.com/@samim/obama-rnn-machine-generated-political-speeches-c8abd18a2ea0

#### **Generating text with an RNN Language Model**

#### Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on Harry Potter:



"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.

Source: https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6

#### **Generating text with an RNN Language Model**

#### Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on recipes:

Title: CHOCOLATE RANCH BARBECUE
Categories: Game, Casseroles, Cookies, Cookies
Yield: 6 Servings



- 2 tb Parmesan cheese -- chopped
- 1 c Coconut milk
- 3 Eggs, beaten

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.

Source: https://gist.github.com/nylki/1efbaa36635956d35bcc

#### **Generating text with a RNN Language Model**

#### Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on paint color names:



This is an example of a character-level RNN-LM (predicts what character comes next)

Source: http://aiweirdness.com/post/160776374467/new-paint-colors-invented-by-neural-network

#### **Evaluating Language Models**

The standard evaluation metric for Language Models is perplexity.

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)}|\ \boldsymbol{x}^{(t)},\dots,\boldsymbol{x}^{(1)})} \right)^{1/T}$$
 Normalized by number of word

• This is equal to the exponential of the cross-entropy loss  $J(\theta)$ :

$$= \prod_{t=1}^T \left(\frac{1}{\hat{y}_{x_{t+1}}^{(t)}}\right)^{1/T} = \exp\left(\frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{x_{t+1}}^{(t)}\right) = \exp(J(\theta))$$

**Lower** perplexity is better!

# RNNs greatly improved perplexity over what came before

	Model	Perplexity
Increasingly complex RNNs	Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
	RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
	RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
	Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
	LSTM-2048 (Jozefowicz et al., 2016)	43.7
	2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
	Ours small (LSTM-2048)	43.9
	Ours large (2-layer LSTM-2048)	39.8

Perplexity improves (lower is better)

Source: https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/

#### 5. Recap

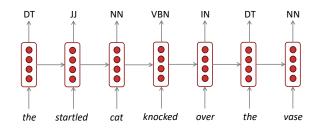
- Language Model: A system that predicts the next word
- Recurrent Neural Network: A family of neural networks that:
  - Take sequential input of any length
  - Apply the same weights on each step
  - Can optionally produce output on each step
- Recurrent Neural Network ≠ Language Model
- We've shown that RNNs are a great way to build a LM (despite some problems)
- RNNs are also useful for much more!

#### Why should we care about Language Modeling?

- Language Modeling is a benchmark task that helps us measure our progress on predicting language use
- Language Modeling is a subcomponent of many NLP tasks, especially those involving generating text or estimating the probability of text:
  - · Predictive typing
  - Speech recognition
  - · Handwriting recognition
  - Spelling/grammar correction
  - · Authorship identification
  - · Machine translation
  - Summarization
  - Dialogue
  - . . . .
  - · etc.

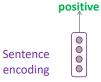
Everything else in NLP has now been rebuilt upon Language Modeling: GPT-3 is an LM!

# Other RNN uses: RNNs can be used for sequence tagging e.g., part-of-speech tagging, named entity recognition

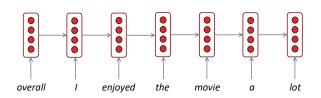


#### RNNs can be used for sentence classification

#### e.g., sentiment classification

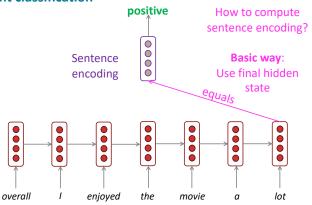


How to compute sentence encoding?

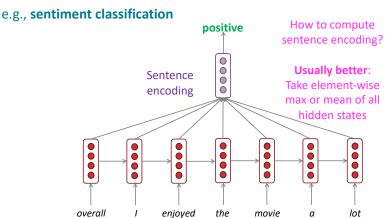


#### RNNs can be used for sentence classification

e.g., sentiment classification

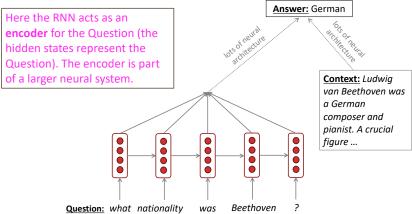


#### RNNs can be used for sentence classification



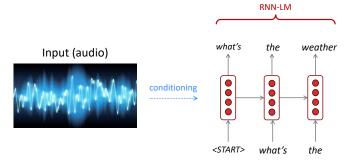
#### RNNs can be used as an encoder module

e.g., **question answering**, machine translation, *many other tasks!* 



#### RNN-LMs can be used to generate text

e.g., speech recognition, machine translation, summarization



This is an example of a *conditional language model*. We'll see Machine Translation in much more detail starting next lecture.

## Terminology and a look forward

The RNN described in this lecture = simple/vanilla/Elman RNN



Next lecture: You will learn about other RNN flavors

like LSTM

and GRU



and multi-layer RNNs



**By the end of the course:** You will understand phrases like "stacked bidirectional LSTMs with residual connections and self-attention"



# Summary of Language Modelling and RNNs

- Modelling the distribution over strings of words in natural language is a fundamental task
- This task can be decomposed into predicting the next word given the prefix of previous words
- Recurrent neural networks (RNNs) can encode arbitrarily long prefixes by passing information through arbitrarily many hidden representations
- Log-likelihood training gives probability estimates for the next word
- Using the same weights for every position allows RNNs to generalise across sequence positions and lengths

## **Outline**

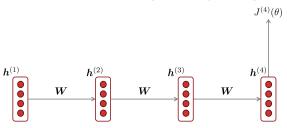
Language Modelling

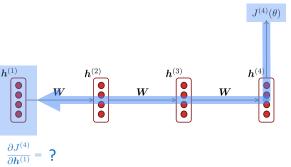
Recurrent Neural Networks

Deep Backprop with RNNs (LSTMs)

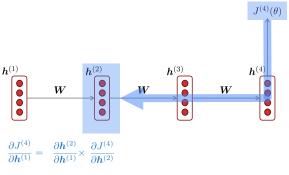
Bigger is Better with RNNs

## 1. Problems with RNNs: Vanishing and Exploding Gradients

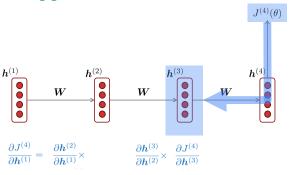




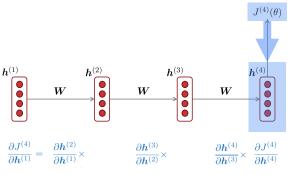
5



chain rule!



chain rule!



chain rule!

#### Vanishing gradient intuition $J^{(4)}(\theta)$ $h^{(3)}$ $h^{(1)}$ $h^{(2)}$ $h^{(4)}$ W WW $\partial J^{(4)}$ $\partial h^{(2)}$ $\partial h^{(3)}$ $\partial J^{(4)}$ $\partial h^{(4)}$ $\overline{\partial m{h}^{(1)}}$ $\overline{\partial h^{(1)}}$ $\overline{\partial h^{(2)}}$ $\overline{\partial h^{(4)}}$ $\partial h^{(3)}$ Vanishing gradient problem: When these are small, the gradient

What happens if these are small?

Vanishing gradient problem: When these are small, the gradient signal gets smaller and smaller as it backpropagates further

#### Vanishing gradient proof sketch (linear case)

Recall:

$$\boldsymbol{h}^{(t)} = \sigma \left( \boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1 \right)$$

• What if  $\sigma$  were the identity function,  $\sigma(x) = x$  ?

$$rac{\partial m{h}^{(t)}}{\partial m{h}^{(t-1)}} = \mathrm{diag}\left(\sigma'\left(m{W}_hm{h}^{(t-1)} + m{W}_xm{x}^{(t)} + m{b}_1
ight)
ight)m{W}_h \qquad \qquad ext{(chain rule)}$$

$$= m{I} \ m{W}_h = \ m{W}_h$$

• Consider the gradient of the loss  $J^{(i)}(\theta)$  on step i, with respect to the hidden state  ${m h}^{(j)}$  on some previous step j. Let  $\ell=i-j$ 

$$\begin{split} \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} & \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \leq i} \boldsymbol{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \boldsymbol{\underline{W}}_h^{\ell} & \text{(value of } \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \text{)} \end{split}$$

If  $W_h$  is "small", then this term gets exponentially problematic as  $\ell$  becomes large

Source: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. <a href="https://proceedings.mlr.press/v28/pascanu13.pdf">https://proceedings.mlr.press/v28/pascanu13.pdf</a> (and supplemental materials), at <a href="https://proceedings.mlr.press/v28/pascanu13-supp.pdf">https://proceedings.mlr.press/v28/pascanu13-supp.pdf</a>

# Vanishing gradient proof sketch (linear case)

• What's wrong with  $W_h^{\ell}$ ?

- sufficient but not necessary
- Consider if the eigenvalues of  $W_h$  are all less than 1:

$$\lambda_1,\lambda_2,\dots,\lambda_n < 1$$
  $q_1,q_2,\dots,q_n$  (eigenvectors)

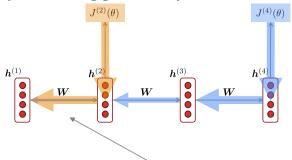
• We can write  $\frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} W_h^{\ell}$  using the eigenvectors of  $W_h$  as a basis:

$$\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}}.\boldsymbol{W}_{h}^{\ell} = \sum_{i=1}^{n} c_{i} \overline{\lambda_{i}^{\ell}} \boldsymbol{q}_{i} \approx \boldsymbol{0} \text{ (for large $\ell$)}$$

Approaches 0 as  $\ell$  grows, so gradient vanishes

- What about nonlinear activations σ (i.e., what we use?)
  - Pretty much the same thing, except the proof requires  $\lambda_i < \gamma$  for some  $\gamma$  dependent on dimensionality and  $\sigma$

#### Why is vanishing gradient a problem?



Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are basically updated only with respect to near effects, not long-term effects.

#### **Effect of vanishing gradient on RNN-LM**

- LM task: When she tried to print her tickets, she found that the printer was out of toner.
   She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her \_\_\_\_\_\_
- To learn from this training example, the RNN-LM needs to model the dependency between "tickets" on the 7<sup>th</sup> step and the target word "tickets" at the end.
- But if the gradient is small, the model can't learn this dependency
  - So, the model is unable to predict similar long-distance dependencies at test time
- In practice a simple RNN will only condition ~7 tokens back [vague rule-of-thumb]

#### Why is exploding gradient a problem?

If the gradient becomes too big, then the SGD update step becomes too big:

$$heta^{new} = heta^{old} - \overbrace{lpha}^{ ext{learning rate}} \int_{ ext{gradient}}^{ ext{gradient}}^{ ext{learning rate}}$$

- This can cause bad updates: we take too large a step and reach a weird and bad parameter configuration (with large loss)
  - You think you've found a hill to climb, but suddenly you're in Iowa
- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)

#### **Gradient clipping: solution for exploding gradient**

 Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

```
Algorithm 1 Pseudo-code for norm clipping  \hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta} 
 \mathbf{if} \quad \|\hat{\mathbf{g}}\| \geq threshold \ \mathbf{then} 
 \hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}} 
 \mathbf{end} \quad \mathbf{if}
```

- Intuition: take a step in the same direction, but a smaller step
- In practice, remembering to clip gradients is important, but exploding gradients are an easy problem to solve

Source: "On the difficulty of training recurrent neural networks", Pascanu et al. 2013, http://proceedings.mlr.press/v28/pascanu13.pdf

#### How to fix the vanishing gradient problem?

- The main problem is that it's too difficult for the RNN to learn to preserve information over many timesteps.
- In a vanilla RNN, the hidden state is constantly being rewritten

$$oldsymbol{h}^{(t)} = \sigma \left( oldsymbol{W}_h oldsymbol{h}^{(t-1)} + oldsymbol{W}_x oldsymbol{x}^{(t)} + oldsymbol{b} 
ight)$$

Could we design an RNN with separate memory which is added to?

#### Long Short-Term Memory RNNs (LSTMs)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the problem of vanishing gradients
  - Everyone cites that paper but really a crucial part of the modern LSTM is from Gers et al. (2000)



- Only started to be recognized as promising through the work of S's student Alex Graves c. 2006
  - Work in which he also invented CTC (connectionist temporal classification) for speech recognition
- But only really became well-known after Hinton brought it to Google in 2013
  - Following Graves having been a postdoc with Hinton

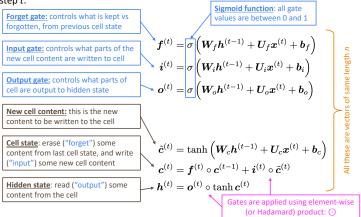
Hochreiter and Schmidhuber, 1997. Long short-term memory. https://www.bioinf.jku.at/publications/older/2604.pdf Gers, Schmidhuber, and Cummins, 2000. Learning to Forget: Continual Prediction with LSTM. https://dl.acm.org/doi/10.1162/089976600300015015 Graves, Fernandez, Gomez, and Schmidhuber, 2006. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural nets. https://www.cs.toronto.edu/~grayes/icml 2006.pdf

#### Long Short-Term Memory RNNs (LSTMs)

- On step t, there is a hidden state  $m{h}^{(t)}$  and a cell state  $m{c}^{(t)}$ 
  - Both are vectors length n
  - The cell stores long-term information
  - The LSTM can read, erase, and write information from the cell
    - · The cell becomes conceptually rather like RAM in a computer
- The selection of which information is erased/written/read is controlled by three corresponding gates
  - The gates are also vectors of length n
  - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between
  - The gates are dynamic: their value is computed based on the current context

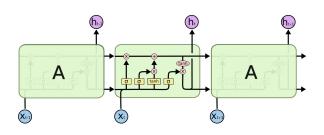
# Long Short-Term Memory (LSTM)

We have a sequence of inputs  $x^{(t)}$ , and we will compute a sequence of hidden states  $h^{(t)}$  and cell states  $c^{(t)}$ . On timestep t:



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

You can think of the LSTM equations visually like this:

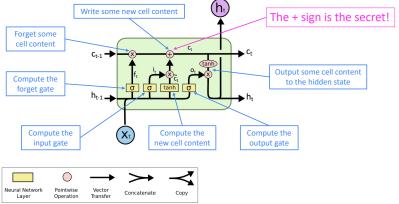




Source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/

## Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



### How does LSTM solve vanishing gradients?

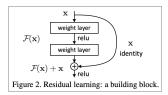
- The LSTM architecture makes it much easier for an RNN to preserve information over many timesteps
  - e.g., if the forget gate is set to 1 for a cell dimension and the input gate set to 0, then the information of that cell is preserved indefinitely.
  - In contrast, it's harder for a vanilla RNN to learn a recurrent weight matrix W<sub>b</sub> that preserves info in the hidden state
  - In practice, you get about 100 timesteps rather than about 7
- However, there are alternative ways of creating more direct and linear pass-through connections in models for long distance dependencies

## Is vanishing/exploding gradient just an RNN problem?

- No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially very deep ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus, lower layers are learned very slowly (i.e., are hard to train)
- Another solution: lots of new deep feedforward/convolutional architectures add more direct connections (thus allowing the gradient to flow)

#### For example:

- · Residual connections aka "ResNet"
- Also known as skip-connections
- The identity connection preserves information by default
- This makes deep networks much easier to train

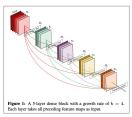


"Deep Residual Learning for Image Recognition", He et al, 2015. https://arxiv.org/pdf/1512.03385.pdf

## Is vanishing/exploding gradient just a RNN problem?

#### Other methods:

- Dense connections aka "DenseNet"
- Directly connect each layer to all future layers!



- Highway connections aka "HighwayNet"
- Similar to residual connections, but the identity connection vs the transformation layer is controlled by a dynamic gate
- Inspired by LSTMs, but applied to deep feedforward/convolutional networks



Conclusion: Though vanishing/exploding gradients are a general problem, RNNs are particularly unstable
due to the repeated multiplication by the same weight matrix [Bengio et al, 1994]

"Densely Connected Convolutional Networks", Huang et al, 2017. https://arxiv.org/pdf/1608.06993.pdf

"Highway Networks", Srivastava et al, 2015. https://arxiv.org/pdf/1505.00387.pdf

"Learning Long-Term Dependencies with Gradient Descent is Difficult", Bengio et al. 1994, http://ai.dinfo.unifi.it/paolo//ps/tnn-94-gradient.pdf

#### LSTMs: real-world success

- In 2013–2015, LSTMs started achieving state-of-the-art results
  - Successful tasks include handwriting recognition, speech recognition, machine translation, parsing, and image captioning, as well as language models
  - LSTMs became the dominant approach for most NLP tasks
- Now (2019–2023), Transformers have become dominant for all tasks
  - For example, in **WMT** (a Machine Translation conference + competition):
    - In WMT 2014, there were 0 neural machine translation systems (!)
    - In WMT 2016, the summary report contains "RNN" 44 times (and these systems won)
    - In WMT 2019: "RNN" 7 times. "Transformer" 105 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, http://www.statmt.org/wmt16/pdf/W15-2301.pdf
Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, http://www.statmt.org/wmt18/pdf/WMT028.ddf
Source: "Findings of the 2019 Conference on Machine Translation (WMT19)", Barrault et al. 2019, http://www.statmt.org/wmt18/pdf/WMT028.ddf
Source: "Findings of the 2019 Conference on Machine Translation (WMT19)", Barrault et al. 2019, http://www.statmt.org/wmt18/pdf/WMT028.ddf

# Summary of Deep Backprop with RNNs

- Backprop through time can be arbitrarily deep, which can cause vanishing or exploding gradients
- Gated identity recurrent connections can address vanishing gradients
- Gradient clipping can address exploding gradients
- LSTMs are effective ways to learn (relatively) long dependencies
- Similar problems exist in all deep neural networks, requiring skip connections or identity biases

## **Outline**

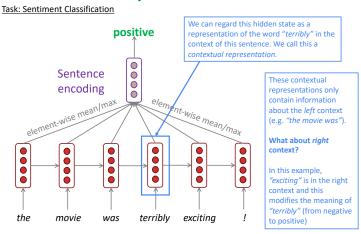
Language Modelling

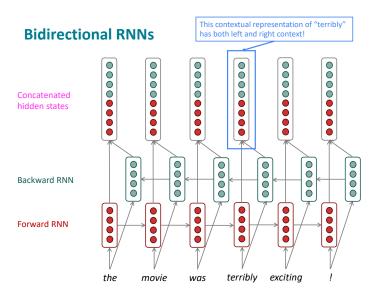
Recurrent Neural Networks

Deep Backprop with RNNs (LSTMs)

Bigger is Better with RNNs

### 4. Bidirectional and Multi-layer RNNs: motivation





#### **Bidirectional RNNs**

On timestep t:

This is a general notation to mean "compute one forward step of the RNN" – it could be a simple RNN or LSTM computation.

Forward RNN 
$$\overrightarrow{\boldsymbol{h}}^{(t)} = \overline{\text{RNN}_{\text{FW}}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$$
Backward RNN  $\overleftarrow{\boldsymbol{h}}^{(t)} = \overline{\text{RNN}_{\text{BW}}}(\overleftarrow{\boldsymbol{h}}^{(t+1)}, \boldsymbol{x}^{(t)})$ 
Generally, these two RNNs have separate weights

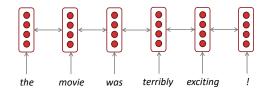
Concatenated hidden states  $\boxed{\boldsymbol{h}^{(t)}} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$ 

We regard this as "the hidden state" of a bidirectional RNN.

This is what we pass on to the next parts of the network.

34

## **Bidirectional RNNs: simplified diagram**



The two-way arrows indicate bidirectionality and the depicted hidden states are assumed to be the concatenated forwards+backwards states

#### **Bidirectional RNNs**

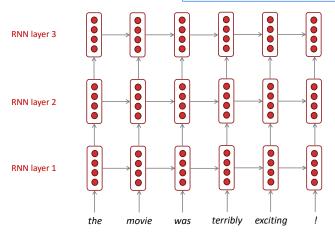
- Note: bidirectional RNNs are only applicable if you have access to the entire input sequence
  - They are not applicable to Language Modeling, because in LM you only have left context available.
- If you do have entire input sequence (e.g., any kind of encoding), bidirectionality is powerful (you should use it by default).
- For example, BERT (Bidirectional Encoder Representations from Transformers) is a
  powerful pretrained contextual representation system built on bidirectionality.
  - You will learn more about transformers, including BERT, in a couple of weeks!

### **Multi-layer RNNs**

- RNNs are already "deep" on one dimension (they unroll over many timesteps)
- We can also make them "deep" in another dimension by applying multiple RNNs – this is a multi-layer RNN.
- This allows the network to compute more complex representations
  - The lower RNNs should compute lower-level features and the higher RNNs should compute higher-level features.
- Multi-layer RNNs are also called stacked RNNs.

## **Multi-layer RNNs**

The hidden states from RNN layer *i* are the inputs to RNN layer *i*+1



#### **Multi-layer RNNs in practice**

- Multi-layer or stacked RNNs allow a network to compute more complex representations
  - they work better than just have one layer of high-dimensional encodings!
  - The lower RNNs should compute lower-level features and the higher RNNs should compute higher-level features.
- High-performing RNNs are usually multi-layer (but aren't as deep as convolutional or feed-forward networks)
- For example: In a 2017 paper, Britz et al. find that for Neural Machine Translation, 2 to 4 layers is best for the encoder RNN, and 4 layers is best for the decoder RNN
  - Often 2 layers is a lot better than 1, and 3 might be a little better than 2
  - Usually, skip-connections/dense-connections are needed to train deeper RNNs (e.g., 8 layers)
- Transformer-based networks (e.g., BERT) are usually deeper, like 12 or 24 layers.
  - You will learn about Transformers later; they have a lot of skipping-like connections

# Summary of Bigger is Better with RNNs

- Bidirectional LSTMs (BiLSTMs) are useful as sentence encoders
- Stacking RNNs adds depth, giving more complex abstract features