



Natural Language Processing

Lecture 06 Language Models; N-gram LMs; Recurrent Neural LMs

Qun Liu, Valentin Malykh
Huawei Noah's Ark Lab



Spring 2020
A course delivered at MIPT, Moscow



Content

- 1 Language Models (LMs)
- 2 N-gram Language Models
- 3 Language Models Based on Recurrent Neural Networks
- 4 Neural Networks Tips and Tricks

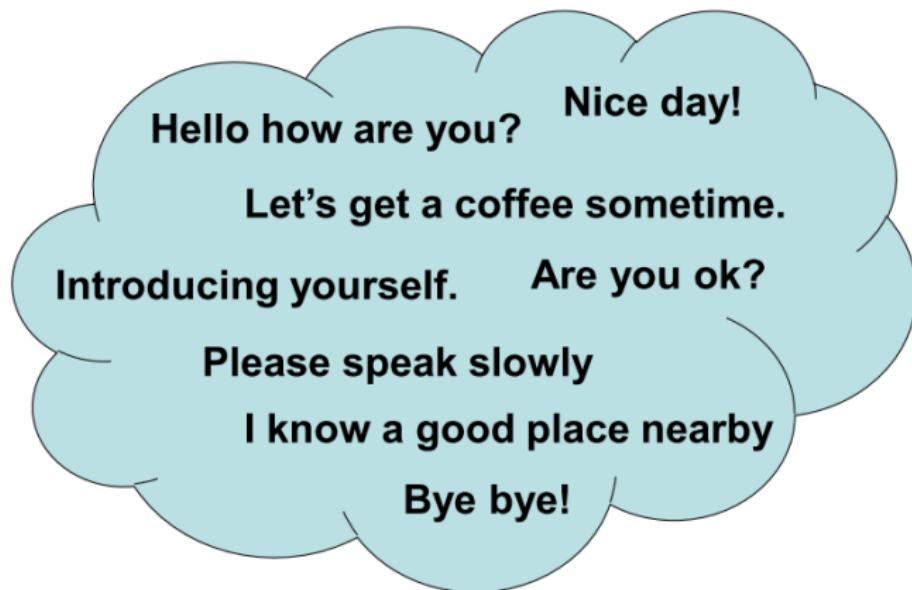


Content

- 1 Language Models (LMs)
- 2 N-gram Language Models
- 3 Language Models Based on Recurrent Neural Networks
- 4 Neural Networks Tips and Tricks



How can we define a language? (Recap)





How can we define a language? (Recap)

— The set theory approach

- A language can be defined as the set of sentences which can be accepted by the speakers of that language.
- It is not possible to define a natural language by enumerate all the sentences, because the number of sentences in a natural languages is infinite.
- Two feasible ways to define a language with infinite sentences:
 - By a Grammar
 - By an Automaton



How can we define a language?

— The probabilistic approach

- A language can also be defined as a probabilistic distribution over all the possible sentences.
- A statistical language model is a probability distribution over sequences of words (sentences) in a given language L :

$$\sum_{s \in V^+} P_{LM}(s) = 1$$

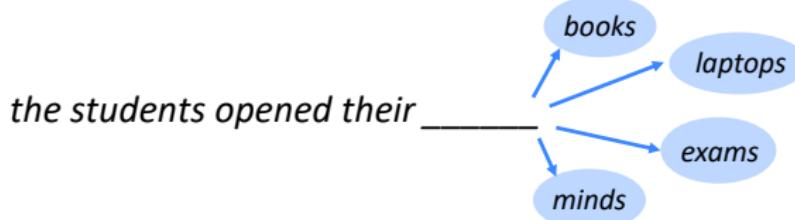
- Or:

$$\sum_{\substack{s=w_1 w_2 \dots w_n \\ w_i \in V, n > 0}} P_{LM}(s) = 1$$



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 $x^{(t+1)}$:

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

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$

- 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 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:

$$P(x^{(1)}, \dots, x^{(T)}) = P(x^{(1)}) \times P(x^{(2)} | x^{(1)}) \times \dots \times P(x^{(T)} | x^{(T-1)}, \dots, x^{(1)})$$

$$= \prod_{t=1}^T P(x^{(t)} | x^{(t-1)}, \dots, x^{(1)})$$



This is what our LM provides



LMs are extremely powerful tools

- The definition of LMs is simple but LMs are extremely powerful tools in NLP
- A number of NLP problems could be viewed as variations of language modelling problems, and can be solved by LMs.
 - Text generation / Data to text / Image captioning / Text summarization
 - Machine translation
 - Speech Recognition
 - Reading Comprehension
 - POS tagging / Entity recognition / Parsing
- Any task which could be transferred to a sequential problems is suitable for LMs



The bag-of-word problem

From a collection of 100 sentences, we considered the 38 sentences with fewer than 11 words each by simply using an n-gram LM.

Exact reconstruction (24 of 38)

- Please give me your response as soon as possible.
⇒ Please give me your response as soon as possible.

Reconstruction preserving meaning (8 of 38)

- Now let me mention some of the disadvantages.
⇒ Let me mention some of the disadvantages now.

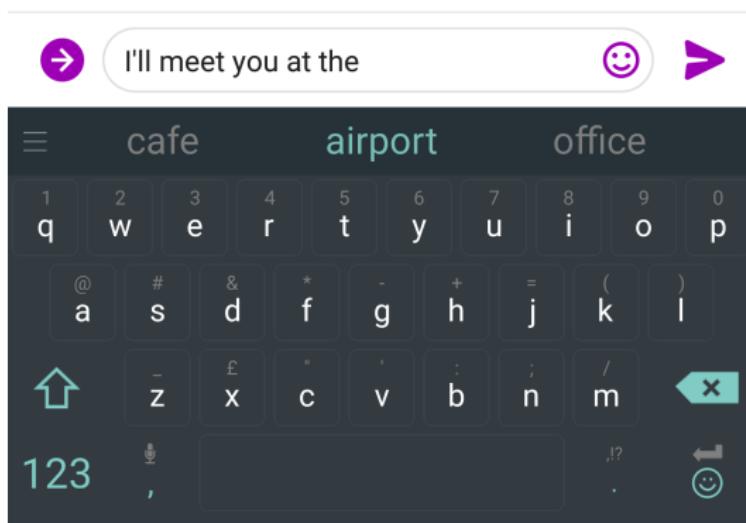
Garbage reconstruction (6 of 38)

- In our organization research has two missions.
⇒ In our missions research organization has two.

Brown P F, Cocke J, Pietra S A D, et al., A statistical approach to machine translation. Computational linguistics, 1990, 16(2):79-85.

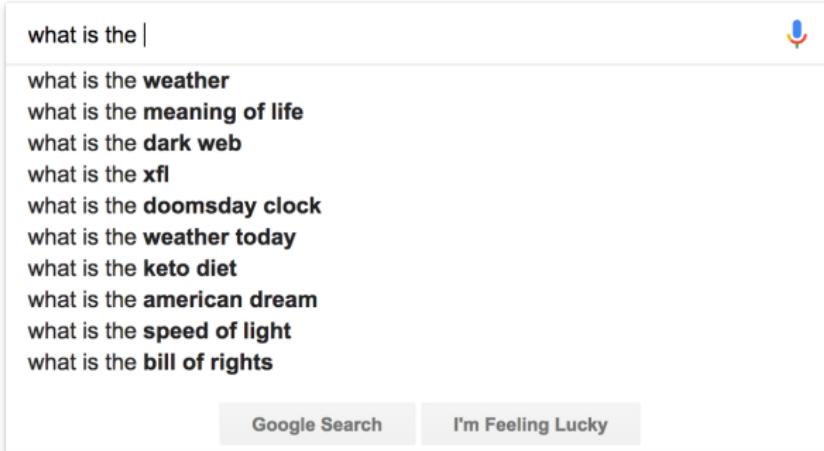


You use Language Models every day!





You use Language Models every day!

The Google logo is displayed in its signature multi-colored font: blue, red, yellow, and green.A screenshot of a Google search interface. The search bar contains the partial query "what is the |". Below the bar is a list of ten suggested completions, each preceded by a small blue square icon. At the bottom are two buttons: "Google Search" and "I'm Feeling Lucky".

- what is the weather
- what is the meaning of life
- what is the dark web
- what is the xfl
- what is the doomsday clock
- what is the weather today
- what is the keto diet
- what is the american dream
- what is the speed of light
- what is the bill of rights





NLP tasks as Conditional LMs

x “input”	w “text output”
An author	A document written by that author
A topic label	An article about that topic
{SPAM, NOT_SPAM}	An email
A sentence in French	Its English translation
A sentence in English	Its French translation
A sentence in English	Its Chinese translation
An image	A text description of the image
A document	Its summary
A document	Its translation
Meteorological measurements	A weather report
Acoustic signal	Transcription of speech
Conversational history + database	Dialogue system response
A question + a document	Its answer
A question + an image	Its answer

Chris Dyer, Conditional LMs (slides)



Content

- 1 Language Models (LMs)
- 2 N-gram Language Models
- 3 Language Models Based on Recurrent Neural Networks
- 4 Neural Networks Tips and Tricks



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:** A *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”
 - 4-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.

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = P(\mathbf{x}^{(t+1)} | \underbrace{\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}}_{n-1 \text{ words}}) \quad (\text{assumption})$$

$$\begin{aligned} \text{prob of a n-gram} &\rightarrow P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \\ \text{prob of a (n-1)-gram} &\rightarrow P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \end{aligned} \quad (\text{definition of conditional prob})$$

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

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})} \quad (\text{statistical approximation})$$



Bigram LM Parameters

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

- Normalize by unigrams:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

- Result:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Dan Jurafsky, Speech and Language Processing (slides)

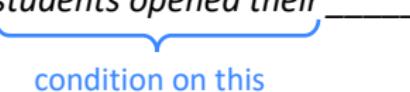


n-gram Language Models: Example

Suppose we are learning a **4-gram** Language Model.

~~as the proctor started the clock, the students opened their _____~~

discard



condition on this

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

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their **books**” occurred 400 times
 - $P(\text{books} | \text{students opened their}) = 0.4$
- “students opened their **exams**” occurred 100 times
 - $P(\text{exams} | \text{students opened their}) = 0.1$

Should we have
discarded the
“proctor” context?



Sparsity Problems with n-gram Language Models

Sparsity Problem 1

Problem: What if “students opened their w ” never occurred in data? Then w has probability 0!

(Partial) Solution: Add small δ to the count for every $w \in V$. This is called *smoothing*.

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

Sparsity Problem 2

Problem: What if “students opened their” never occurred in data? Then we can’t calculate probability for *any* w !

(Partial) Solution: Just condition on “opened their” instead. This is called *backoff*.

Note: Increasing n makes sparsity problems worse.
Typically we can’t have n bigger than 5.



Storage Problems with n-gram Language Models

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

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

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*

today the _____

Business and financial news

get probability distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem:
not much granularity
in the probability
distribution

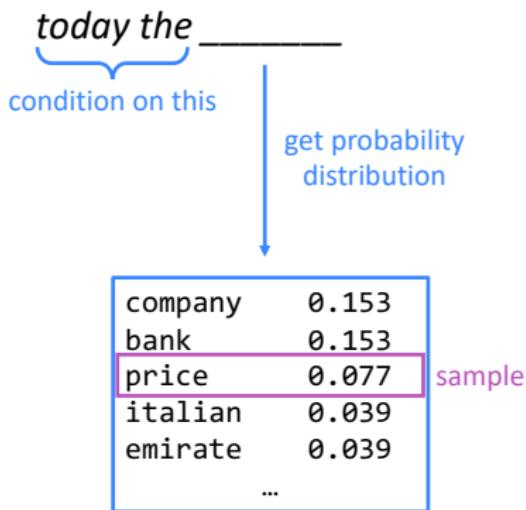
Otherwise, seems reasonable!

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



Generating text with a n-gram Language Model

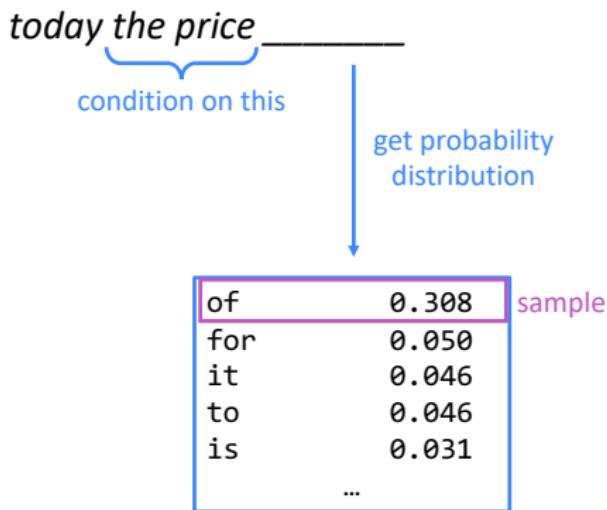
- You can also use a Language Model to generate text.





Generating text with a n-gram Language Model

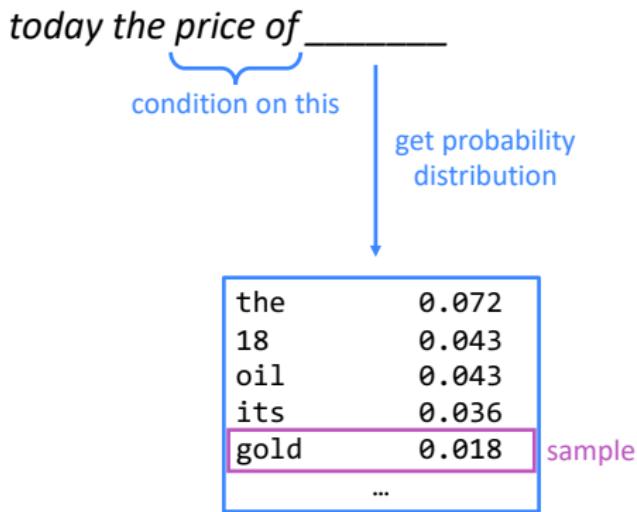
- You can also use a Language Model to generate text.





Generating text with a n-gram Language Model

- You can also use a Language Model to generate text.





Generating text with a n-gram Language Model

- You can also use a Language Model to generate text.

today the price of gold _____



Generating text with a n-gram Language Model

- 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...

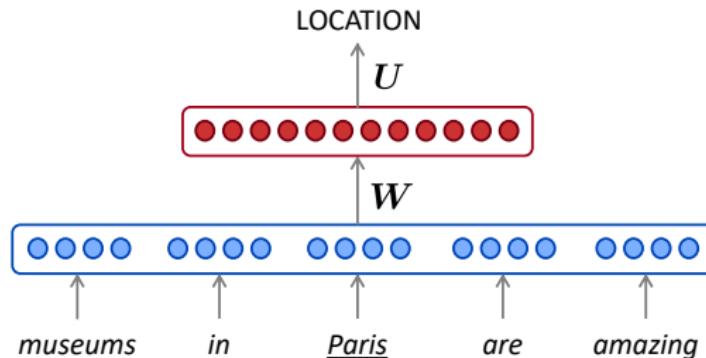


Content

- 1 Language Models (LMs)
- 2 N-gram Language Models
- 3 Language Models Based on Recurrent Neural Networks
- 4 Neural Networks Tips and Tricks

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 $P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$
- How about a **window-based neural model**?
 - We saw this applied to Named Entity Recognition in Lecture 3:





A fixed-window neural Language Model

as the proctor started the clock the students opened their _____
discard 

A fixed-window neural Language Model

output distribution

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

hidden layer

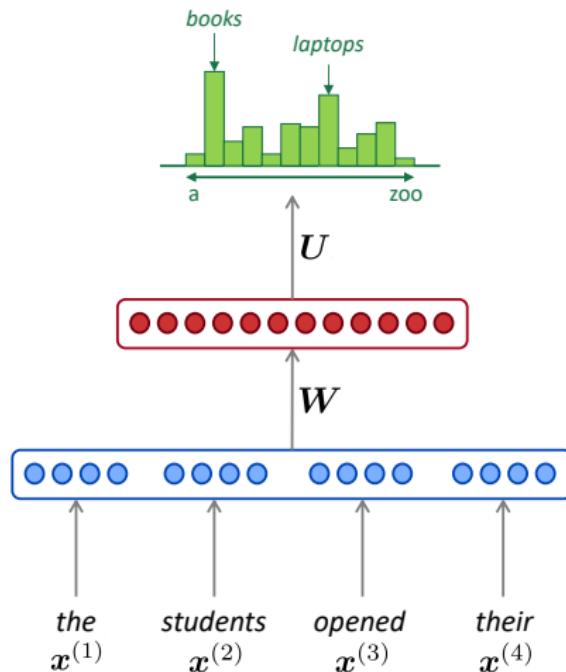
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

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

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$



A fixed-window neural Language Model

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

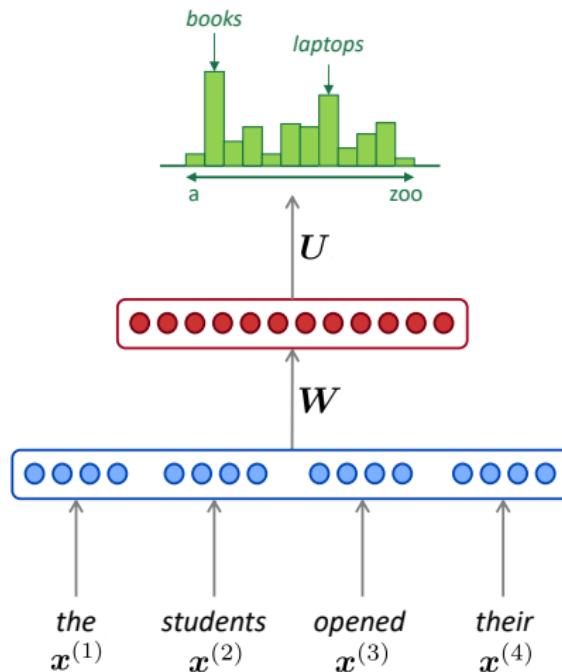
Improvements over n -gram LM:

- 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!
- $x^{(1)}$ and $x^{(2)}$ 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*





Content

3

Language Models Based on Recurrent Neural Networks

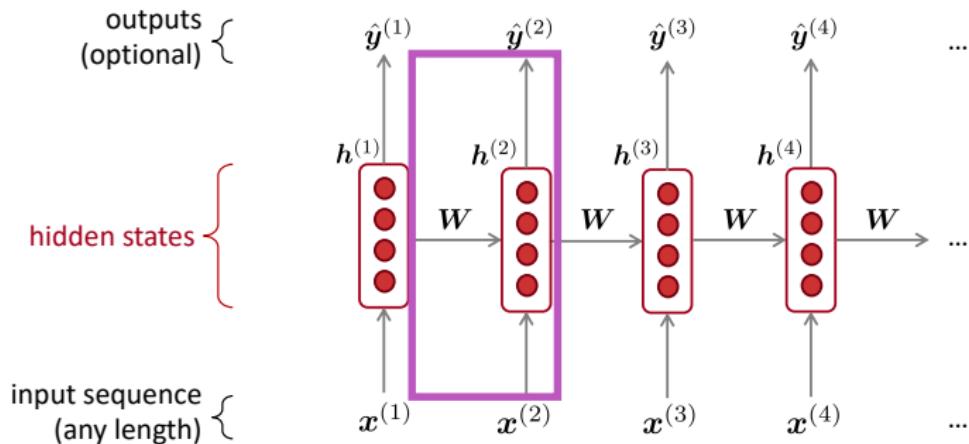
- Recurrent Neural Networks (RNNs)
 - Training a neural language model
 - Other applications of neural language models
 - Evaluation of language models
 - Gradient vanishing and exploding
 - Long Short Term Memory (LSTM)
 - Gated Recurrent Units (GRUs)
 - Vanishing and exploding gradient for other neural networks
 - Bi-directional RNNs and multi-layer RNNs



Recurrent Neural Networks (RNN)

A family of neural architectures

Core idea: Apply the same weights W repeatedly

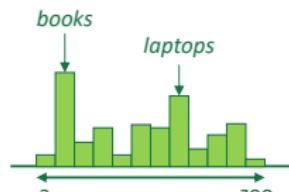




A Simple RNN Language Model

output distribution

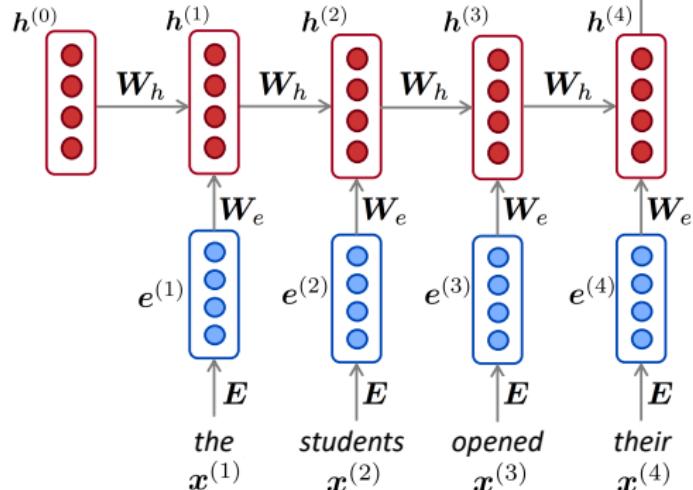
$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$



hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

$\mathbf{h}^{(0)}$ is the initial hidden state



word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$

Note: this input sequence could be much longer, but this slide doesn't have space!

35



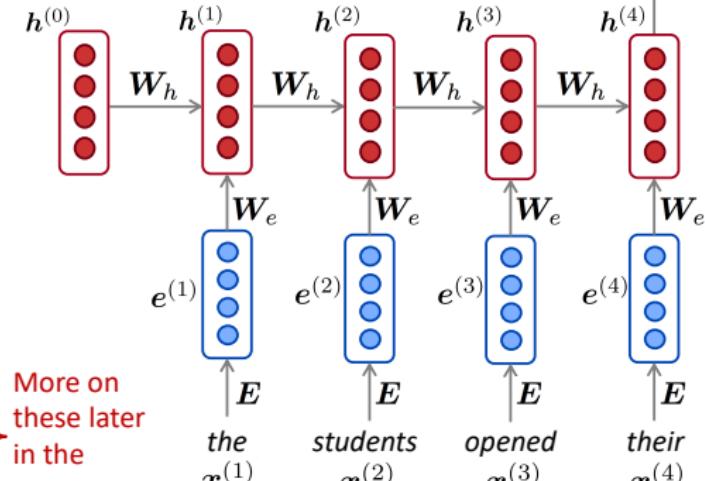
RNN Language Models

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
- 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**



36

Christopher Manning, Natural Language Processing with Deep Learning, Standford U. CS224n





Content

3

Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- **Training a neural language model**
- Other applications of neural language models
- Evaluation of language models
- Gradient vanishing and exploding
- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRUs)
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs



Training an RNN Language Model

- 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 $x^{(t+1)}$):

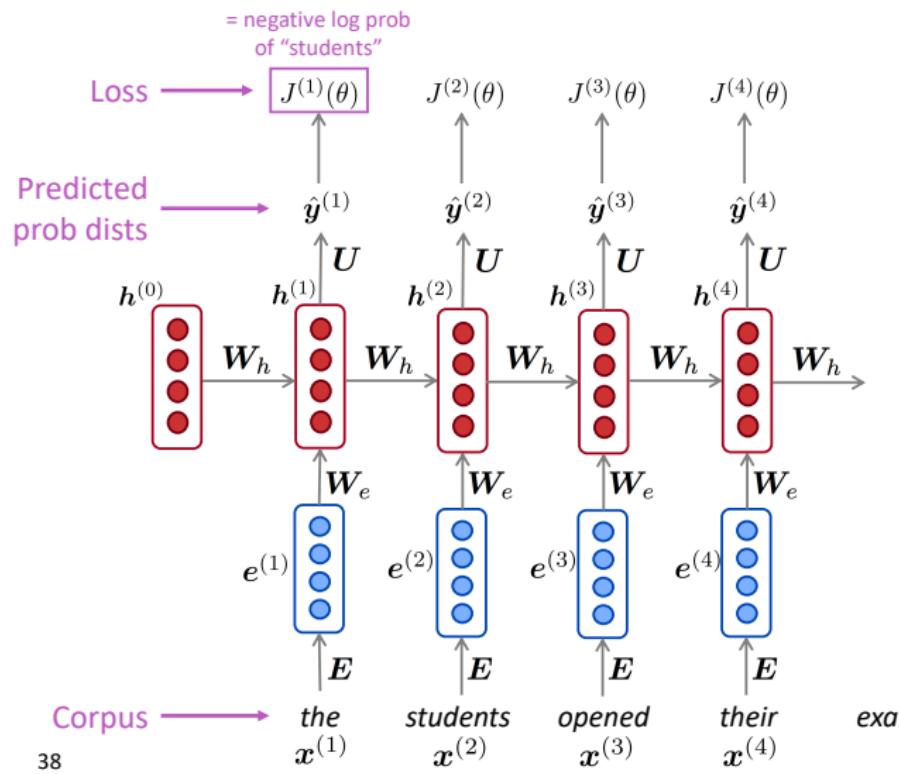
$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$

- Average this to get **overall loss** for entire training set:

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

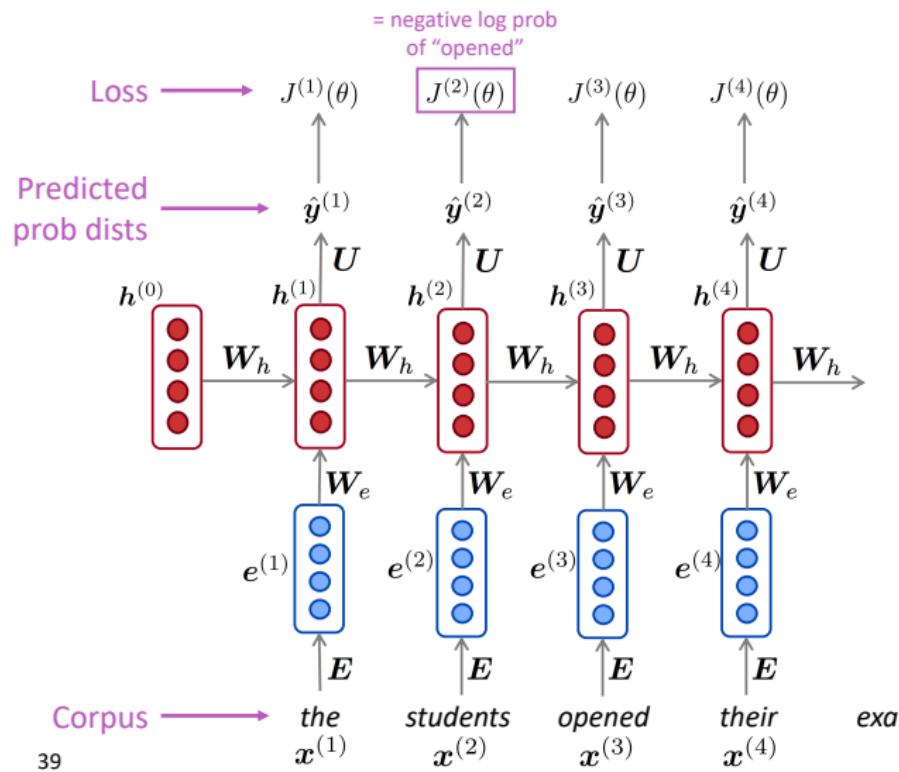


Training an RNN Language Model



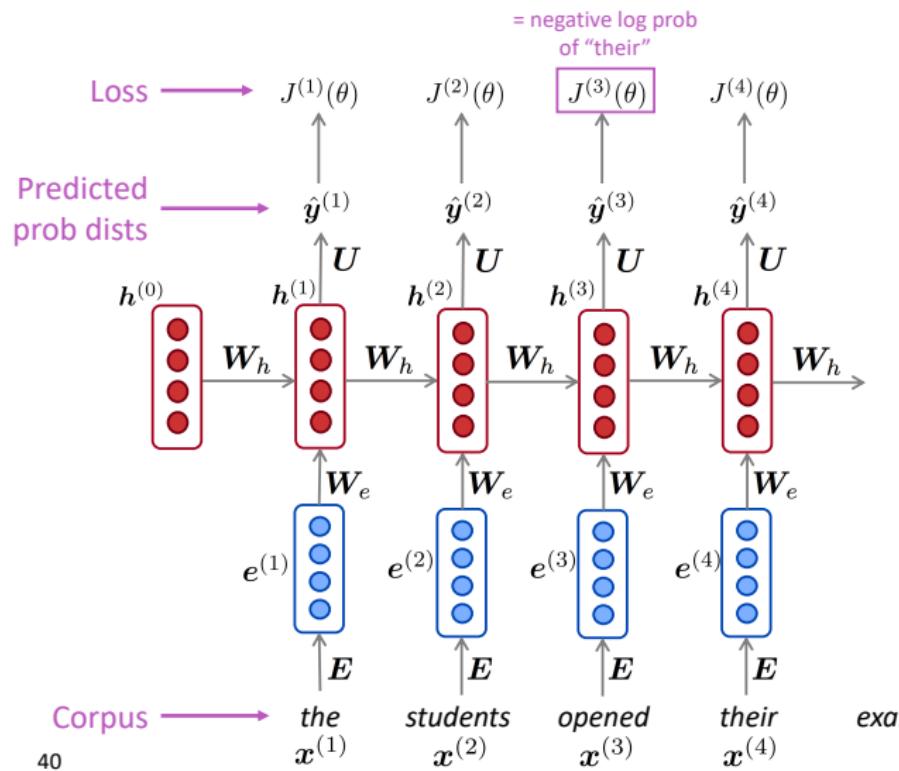


Training an RNN Language Model



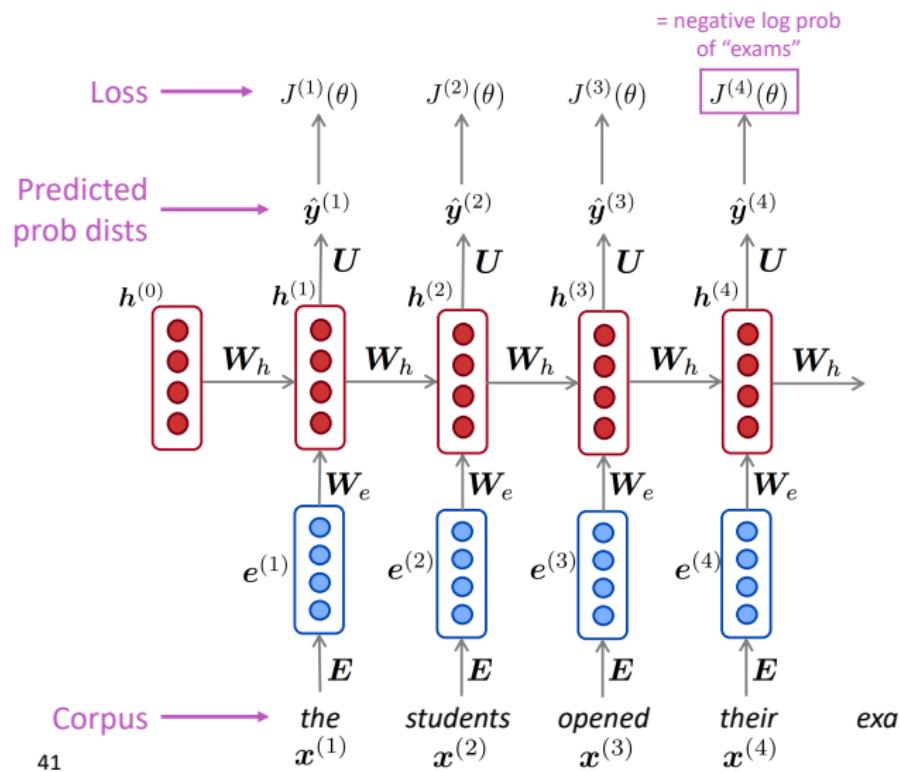


Training an RNN Language Model





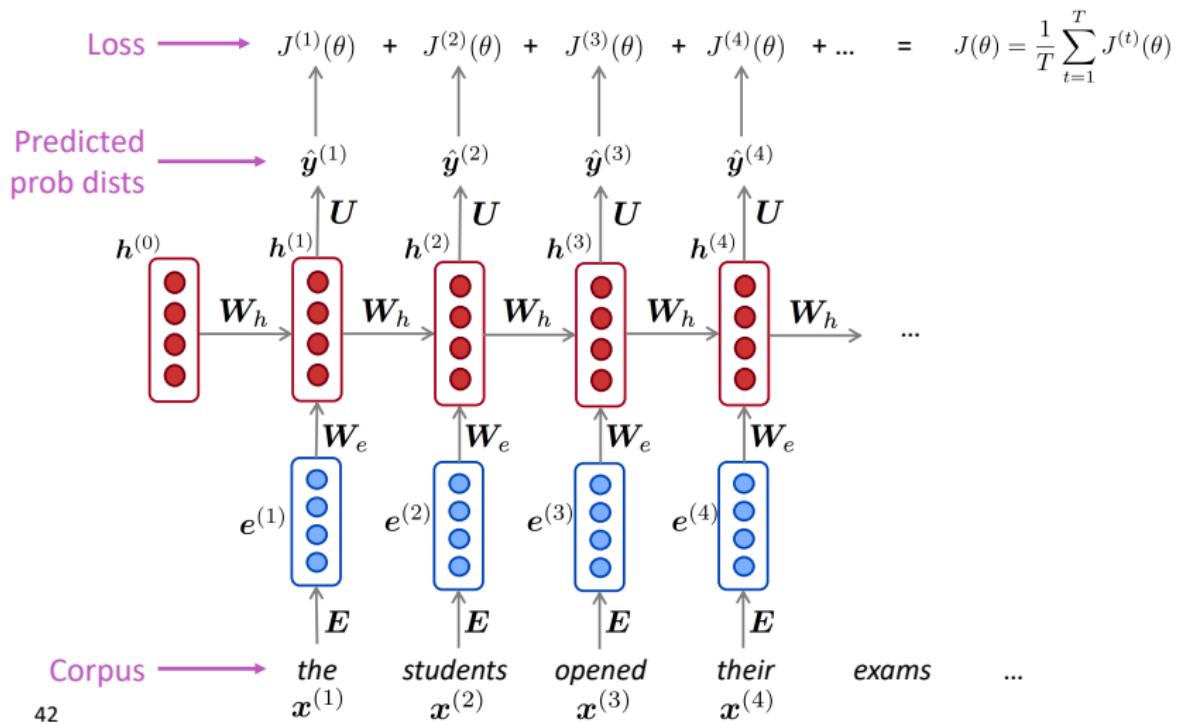
Training an RNN Language Model





Training an RNN Language Model

“Teacher forcing”





Training a RNN Language Model

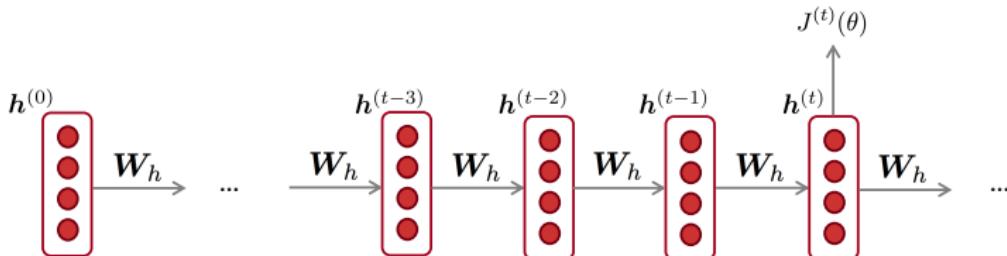
- However: Computing loss and gradients across **entire corpus** $x^{(1)}, \dots, x^{(T)}$ is **too expensive**!

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

- In practice, consider $x^{(1)}, \dots, x^{(T)}$ 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.



Backpropagation for RNNs



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

$$\text{Answer: } \frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)}$$

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

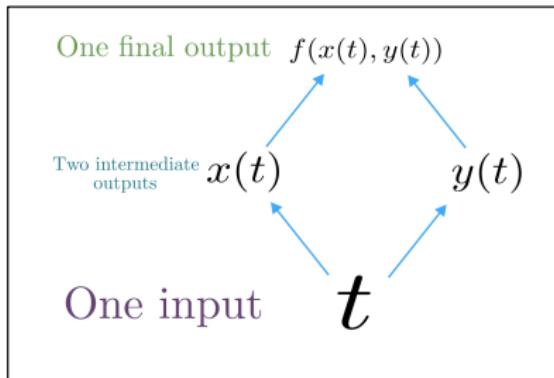
Why?



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:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$



Source:

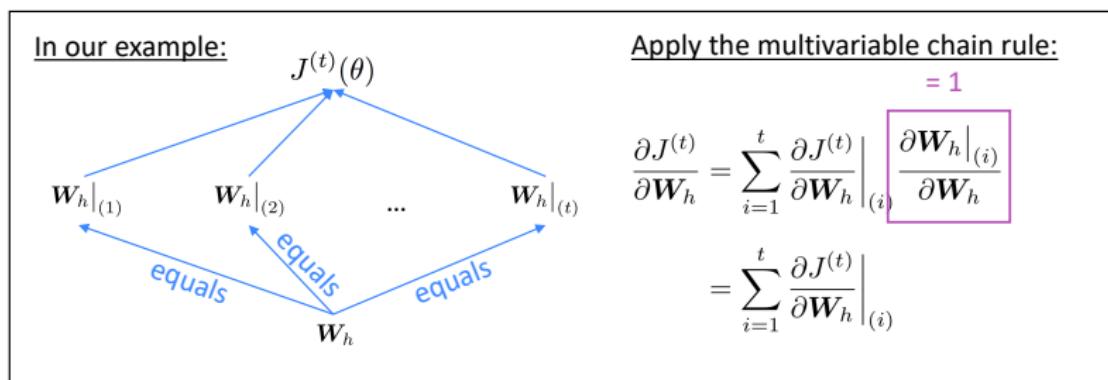
<https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-rule-simple-version>



Backpropagation for RNNs: Proof sketch

- 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:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

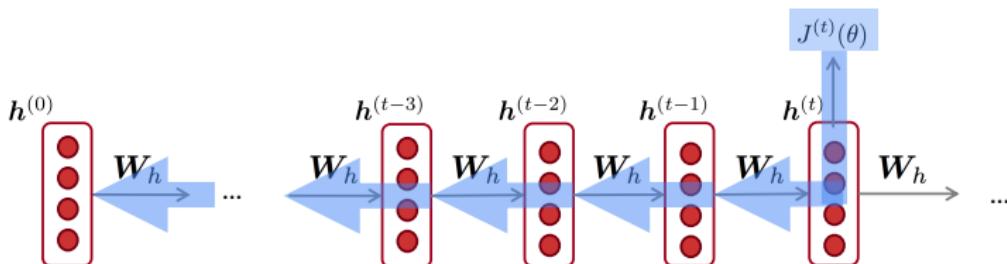


Source:

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



Backpropagation for RNNs



$$\frac{\partial J^{(t)}}{\partial \mathbf{W}_h} = \sum_{i=1}^t \left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right|_{(i)}$$

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

Question: How do we calculate this?



Content

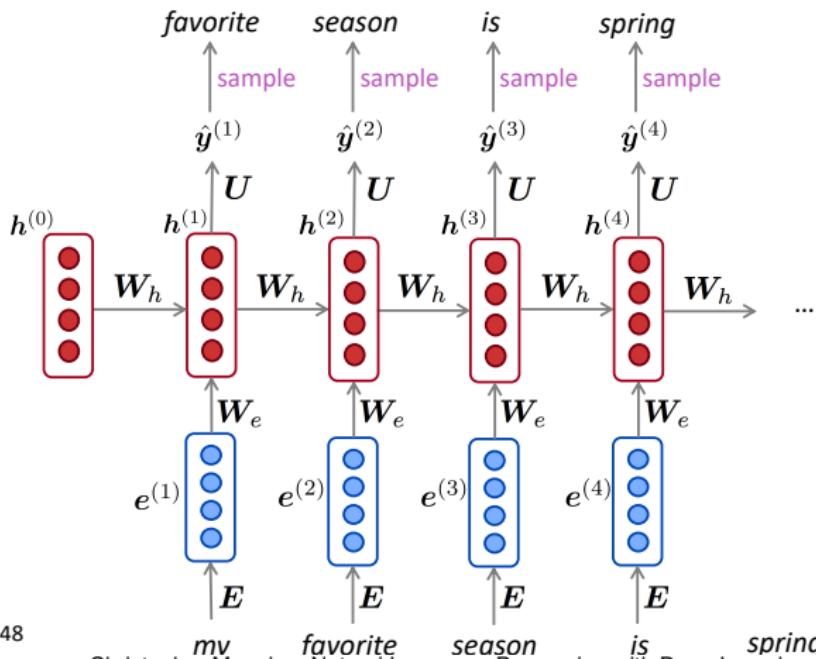
3

Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
 - Evaluation of language models
 - Gradient vanishing and exploding
 - Long Short Term Memory (LSTM)
 - Gated Recurrent Units (GRUs)
 - Vanishing and exploding gradient for other neural networks
 - Bi-directional RNNs and multi-layer RNNs

Generating text with a RNN Language Model

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





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 **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>