Natural Language Processing with Deep Learning CS224N/Ling284



Tatsunori Hashimoto

Lecture 5: Language Models and Recurrent Neural Networks

Lecture Plan

Language modeling + RNNs

• 1. A new NLP task: Language Modeling (20 mins)

motivates

This is the most important concept in the class! Leads to most of modern NLP

- 2. Language models with neural nets: Recurrent Neural Networks (RNNs) (25 mins)
- 4. Problems with RNNs (15 mins) Impor

Important and used in Ass3, but not the only way to build LMs

• 5. Recap on RNNs/LMs (10 mins)

Reminders:

Assignment 2 – Due Jan 23, Thursday

In Assignment 3, out Thursday – Neural MT with RNNs!

1. Language Modeling

Language Modeling is the task of predicting what word comes next



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

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

• 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 $x^{(1)}, \ldots, x^{(T)}$, then the probability of this text (according to the Language Model) is:

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

This is what our LM provides

You use Language Models every day!



You use Language Models every day!



ChatGPT is optimized for dialogue. Our goal is to make AI systems more natural to interact with, and your feedback will help us improve our sy

od.rs X	CHAT COMPOSER
Implement the cleanup function for the tranport stack. Do not make the upgrade listeners optional. Intel Accept Mol Reject Follow-up Instructions ⊕#K	mod.rs Current File Could you make it easier to switch certificates in the transport listeners?
<pre>pub(crate) struct TransportStack { l4: ListemerEndpoint, tls: Option<arc<acceptor>>, // Listemers sent from the old process for graceful upgrade</arc<acceptor></pre>	I'll help modify the code to make certificate switching more flexible. The main changes will be to enhance the TlsAccept trait and modify how certificates are handled in the TlsSettings . Here are the key changes:
#[cfg(unix)]	🐵 mod.rs 🗠 Ask 🗈 Copy 🗅 Apply
upgrade_listeners: Option <listenfds>, upgrade_listeners: ListenFds,</listenfds>	// existing code
<pre>} impl TransportStack { pub fn as_str(&self) -> &str { self.l4.as_str() } pub async fn listen(&mut self) -> Result<()> { self.l4 self</pre>	<pre>/// The APIs to customize things like certificate during TLS ser #[async_trait] pub trait TlsAccept { /// Called during TLS handshake to select and apply certific /// Raturns true if certificates were successfully applied async fn certificate_callback(&Self, SSL Sum tTLSRef) -> be // Default implementation does nothing and indicates no false } /// Called when certificate selection fails</pre>
Some(self.upgrade_listeners.take()),	async fn on_cert_error(&self, _error: &str) {
<pre>} pub async fn accept(&mut self) -> Result<uninitializedstream> { let stream = self.l4.accept().await?; Ok(UninitializedStream { L4: stream, tls: self.tls.clone(), }) } </uninitializedstream></pre>	<pre>// Default imprementation does hotning } // Add a default no-op implementation that can be used when no c #(derive(Default))</pre>

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 been rebuilt upon Language Modeling: ChatGPT is an LM!

What can you do with next-word prediction?

•A sufficiently strong (!) language model can do many, many things

Stanford University is located in _____, California. [Trivia]

I put _____ fork down on the table. [syntax]

The woman walked across the street, checking for traffic over _____ shoulder. [coreference]

I went to the ocean to see the fish, turtles, seals, and _____. [lexical semantics/topic]

Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ____. [sentiment]

Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the _____. [some reasoning – this is harder]

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, ____ [some basic arithmetic]

Next word prediction (2019 – GPT2)

- Fixing /.ssh/ Permissions in Cygwin
- In Cygwin on Windows I found that I could not set the permissions of my /.ssh/ folder to be 0600, as is required for ssh to allow you to use keys. The symptom I had was that I no matter what I did, it always modfied the owner and the group
- , making these files read-only. I tried just starting ssh and using a program like passwd which modifies the owner/group to 0644 permissions, but it didn't work. It was annoying but not nearly as bad as it could be, just annoying to have to fiddle with permissions. So I decided to do this as the reverse: For each .ssh folder in your Cygwin directory, set the permission of each file and directory to 2644. Then start ssh with: ssh-add /.ssh/ The first time you connect, your permissions of the /.ssh folder that you just had in Cygwin will be set to 0600. Then every time you connect again, your permissions will change to 0644.

Stylistically correct, content wise, questionable

Next word prediction (2020 – GPT3)

 $\texttt{Context} \rightarrow \texttt{Title: The_Blitz}$

Background: From the German point of view, March 1941 saw an improvement. The Luftwaffe flew 4,000 sorties that month, including 12 major and three heavy attacks. The electronic war intensified but the Luftwaffe flew major inland missions only on moonlit nights. Ports were easier to find and made better targets. To confuse the British, radio silence was observed until the bombs fell. X- and Y-Gerät beams were placed over false targets and switched only at the last minute. Rapid frequency changes were introduced for X-Gerät, whose wider band of frequencies and greater tactical flexibility ensured it remained effective at a time when British selective jamming was degrading the effectiveness of Y-Gerät.

- Q: How many sorties were flown in March 1941?
- A: 4,000
- Q: When did the Luftwaffe fly inland missions?

A:

Target Completion \rightarrow only on moonlit nights

The model can perform tasks when given a few examples ('in-context learning')

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

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

prob of a n-gram
$$= P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$$
 (definition of orob of a (n-1)-gram
$$= P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$$
 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!

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

n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.



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

For example, suppose that in the corpus:

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

Should we have discarded the "proctor" context?

Sparsity Problems with n-gram Language Models



<u>Note:</u> Increasing *n* makes sparsity problems *worse*. Typically, we can't have *n* bigger than 5.

Storage Problems with n-gram Language Models



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*

Business and financial news today the get probability distribution 0.153 company **Sparsity problem**: bank 0.153 not much granularity price 0.077 in the probability italian 0.039 distribution emirate 0.039

Otherwise, seems reasonable!

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

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...

Evaluating Language Models

• The standard evaluation metric for Language Models is perplexity.

perplexity =
$$\prod_{t=1}^{T} \left(\frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})} \right)^{1/T}$$
 Normalized by number of words Inverse probability of corpus, according to Language Model

• 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!

How to build a *neural* language model?

- Recall the Language Modeling task:
 - Input: sequence of words $oldsymbol{x}^{(1)},oldsymbol{x}^{(2)},\ldots,oldsymbol{x}^{(t)}$
 - Output: prob. dist. of the next word $P(\boldsymbol{x}^{(t+1)} | \ \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})$
- 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



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*



2. Recurrent Neural Networks (RNN)

A family of neural architectures

<u>Core idea:</u> Apply the same weights *W* repeatedly





RNN Language Models

 $oldsymbol{h}^{(0)}$

 \bigcirc

RNN Advantages:

- Can process any length input
- Computation for step t can (in ullettheory) use information from many steps back
- Model size doesn't increase for ٠ 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

- Get a big corpus of text which is a sequence of words $m{x}^{(1)},\ldots,m{x}^{(T)}$
- 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:

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

"Teacher forcing"

 However: Computing loss and gradients across entire corpus x⁽¹⁾,...,x^(T) at once is too expensive (memory-wise)!

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

- In practice, consider $x^{(1)}, \ldots, x^{(T)}$ as a sentence (or a document)
- <u>Recall</u>: 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 ?

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

"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:

$$\underbrace{\frac{d}{dt}f(\boldsymbol{x}(t),\boldsymbol{y}(t))}_{\boldsymbol{d}t} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{d\boldsymbol{x}}{dt} + \frac{\partial f}{\partial \boldsymbol{y}}\frac{d\boldsymbol{y}}{dt}$$

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

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

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

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 *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

3. Problems with RNNs: Vanishing and Exploding Gradients

Vanishing gradient proof sketch (linear case)

Recall:

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

- What if σ were the identity function, $\sigma(x) = x$? $\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = \operatorname{diag} \left(\sigma' \left(W_h h^{(t-1)} + W_x x^{(t)} + b_1 \right) \right) W_h \quad \text{(chain rule)}$ $= I W_h = W_h$
- Consider the gradient of the loss $J^{(i)}(heta)$ on step i, with respect to the hidden state $m{h}^{(j)}$ on some previous step j. Let $\ell=i-j$

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \le i} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} & \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \le i} W_h = \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} W_h^{\ell} & \text{(value of } \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \text{)} \\ & \text{If } W_h \text{ is "small", then this term gets} \end{aligned}$$

exponentially problematic as ℓ becomes large

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

Vanishing gradient proof sketch (linear case)

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

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 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 7th 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

Why is exploding gradient a problem?

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

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$
gradient

- 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

- 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

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

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

- First off next time: How about an RNN with separate memory which is added to?
 LSTMs
- And then: Creating more direct and linear pass-through connections in model
 - Attention, residual connections, etc.

4. 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 still relevant today with the rise of *state space models*