Lecture 24: Attention

NLP Lectures: Part 3 of 4

Harvard IACS

CS109B

Pavlos Protopapas, Mark Glickman, and Chris Tanner



Outline



Outline









Previously, we learned about word embeddings



millions of books

word2vec

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)

Previously, we learned about word embeddings



millions of books

word2vec

word embeddings

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)

"The food was delicious. Amazing!" - 4.8/5
 \$yelp

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)



word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)



word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)

"Daaang. What?! Supa Lit" 📥 4.9/5 💱 yelp



Strengths and weaknesses of word embeddings (type-based)?

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)



Strengths:

- Leverages tons of existing data
- Don't need to depend on our data to create embeddings

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)



lssues:

- Out-of-vocabulary (OOV) words
- Not tailored to this dataset

word embeddings (type-based)

- count-based/DSMs (e.g., SVD, LSA)
- Predictive models (e.g., word2vec, GloVe)

Previously, we learned about word embeddings



contextualized embeddings (token-based)

approaches:

\$\$\$ • Middle Eastern, Tapas/Small Plates, Mediterranean

2. Chalawan 108 Thai, Indonesian, Singaporean



Review #1

4.8

contextualized embeddings (token-based)

approaches:

\$\$\$ • Middle Eastern, Tapas/Small Plates, Mediterranean

2. Chalawan 108 Thai, Indonesian, Singaporean



2.5

contextualized embeddings (token-based)

approaches:

\$\$\$ • Middle Eastern, Tapas/Small Plates, Mediterranean

2. Chalawan 2. Chalawan 108 Thai, Indonesian, Singaporean

Chicken Wings, Chicken Shop



Review #53,781

1.0

contextualized embeddings (token-based)

approaches:

1. Sarma \$\$\$ • Middle Eastern, Tapas/Small Plates, Mediterranean 2. Chalawan Thai, Indonesian, Singaporean Every token in the corpus has a 3. Gustazo Cuban Kitchen & Bar 162 contextualized embedding Cuban, Tapas/Small Plates, Bars 4. Posto * * * * 912 \$\$ • Italian, Pizza 00000 00000 00000 00000 00000 5. Avenue Kitchen + Bar \$\$ • Beer Bar, Pizza, Cocktail Bars 0000 found hair а in 7192. Underdog Hot Chicken 45

Review #53,781

contextualized embeddings (token-based)

approaches:

Chicken Wings, Chicken Shop

• Predictive models (e.g., BiLSTMs, GPT-2, BERT)

the



Review #53,781

contextualized embeddings (token-based)

approaches:



Review #53,781

contextualized embeddings (token-based)

approaches:

Strengths:

- Tailored to your particular corpus

Review #53,781

contextualized embeddings (token-based)

approaches:

Weaknesses:

- May not have enough data to produce good results
- Have to train new model for each use case
- Can't leverage a wealth of existed text data (millions of books)???

Review #53,781

contextualized embeddings (token-based)

approaches:

Weaknesses:

- May not have enough data to produce good results
- Have to train new model for each use case
- Can't leverage a wealth of existed text data (millions of books)???

WRONG! We can leverage millions of books!

Review #53,781

contextualized embeddings (token-based)

approaches:

MOBY-DICK Herman Melvil



Language Modelling

(let's input 1 million documents)





The Free Encyclopedia

Language Modelling

(let's input 1 million documents)



The contextualized embeddings for 1 million docs aren't useful to us for a new task (e.g., predicting Yelp reviews), but the learned weights could be!



Using these "pre-trained" *W* and *V*, we can possibly increase our performance on other tasks (e.g., Yelp reviews), since they're very experienced with producing/capturing "meaning"



RECAP

- Language Modelling may help us for other tasks
- LSTMs do a great job of capturing "*meaning*", which can be used for almost every task
 - Given a sequence of N words, we can produce 1 output
 - Given a sequence of N words, we can produce N outputs

RECAP

- Language Modelling may help us for other tasks
- LSTMs do a great job of capturing "*meaning*", which can be used for almost every task
 - Given a sequence of N words, we can produce 1 output
 - Given a sequence of N words, we can produce N outputs
 - What if we wish to have **M** outputs?

We want to produce a variable-length output (e.g., $n \rightarrow m$ predictions)



Outline









Outline



- If our input is a sentence in Language A, and we wish to translate it to Language B, it is clearly sub-optimal to translate word by word (like our current models are suited to do).
- Instead, let a sequence of tokens be the unit that we ultimately wish to work with (a sequence of length N may emit a sequences of length M)
- Seq2seq models are comprised of **2 RNNs**: 1 encoder, 1 decoder



The final hidden state of the encoder RNN is the initial state of the decoder RNN



The final hidden state of the encoder RNN is the initial state of the decoder RNN





The final hidden state of the encoder RNN is the initial state of the decoder RNN



DECODER RNN



DECODER RNN


DECODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



ENCODER RNN



Training occurs like RNNs typically do; the loss (from the decoder outputs) is calculated, and we update weights all the way to the beginning (encoder)

 $n_{\overline{2}}$

 π_1

The

 $n_{\overline{3}}$

dog

 n_4

ran

Hidden layer

Input layer

ENCODER RNN

brown

DECODER RNN

 \hat{y}_3

 h_3^D

chien

 \hat{y}_4

 h_4^D

brun

 \hat{y}_5

 h_5^D

а

 \hat{y}_6

 h_6^D

couru

 \hat{y}_1

 h_1^D

 $\langle s \rangle$

 \hat{y}_2

 h_2^D

Le



ENCODER RNN

See any issues with this traditional **seq2seq** paradigm?



It's crazy that the entire "meaning" of the 1st sequence is expected to be packed into this one embedding, and that the encoder then never interacts w/ the \hat{y}_1 decoder again. Hands free.

Hidden layer

 \hat{y}_2 \hat{y}_3 \hat{y}_4 \hat{y}_5 \hat{y}_6 h_1^E h_1^D h_2^D h_3^D h_4^D h_2^E h_3^E h_4^E h_5^D h_6^D Input layer chien The brown brun dog Le < s >а couru ran

ENCODER RNN

Instead, what if the decoder, at each step, pays attention to a *distribution* of all of the encoder's hidden states?

Instead, what if the decoder, at each step, pays attention to a *distribution* of all of the encoder's hidden states?

Intuition: when we (humans) translate a sentence, we don't just consume the original sentence then regurgitate in a new language; we continuously look back at the original while focusing on different parts.

Outline



Outline



Q: How do we determine how much to pay attention to each of the encoder's hidden layers?













A: Let's base it on our decoder's current hidden state (our current representation of meaning) and all of the encoder's hidden layers!

Attention (raw scores)

 e_1 1.5

*e*₂ 0.9

*e*₃ 0.2

 $e_4 - 0.5$





A: Let's base it on our decoder's current hidden state (our current representation of meaning) and all of the encoder's hidden layers!

Attention (raw scores)





We multiply each encoder's hidden layer by its a_i^1 attention weights to create a context vector c_1^D

Attention (softmax'd)

 $a_1^1 = 0.51$ $a_2^1 = 0.28$ $a_3^1 = 0.14$ $a_3^1 = 0.07$

REMEMBER: each attention weight a_i^j is based on the decoder's current hidden state, too.



REMEMBER: each attention weight a_i^j is based on the decoder's current hidden state, too.



REMEMBER: each attention weight a_i^j is based on the decoder's current hidden state, too.



REMEMBER: each attention weight a_i^j is based on the decoder's current hidden state, too.



REMEMBER: each attention weight a_i^j is based on the decoder's current hidden state, too.



For convenience, here's the Attention calculation summarized on 1 slide



Photo credit: https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html



Photo credit: https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html


Attention:

- greatly improves seq2seq results
- allows us to visualize the contribution each encoding word gave for each decoder's word



Image source: Fig 3 in <u>Bahdanau et al., 2015</u>



are

ga

Takeaway:

Having a separate encoder and decoder allows for $n \rightarrow m$ length predictions.

Attention is powerful; allows us to conditionally weight our focus

SUMMARY

- LSTMs yielded state-of-the-art results on most NLP tasks (2014-2018)
- seq2seq+Attention was an even more revolutionary idea (Google Translate used it)
- Attention allows us to place appropriate weight to the encoder's hidden states
- But, LSTMs require us to iteratively scan each word and wait until we're at the end before we can do anything

Outline



Outline



Transformer Encoder



Transformer Encoder uses attention on itself (self-attention) to create very rich embeddings which can be used for any task.

BERT is a Bidirectional Transformer Encoder. You can attach a final layer that performs whatever task you're interested in (e.g., Yelp reviews).

Its results are unbelievably good.

BERT (a Transformer variant)

BERT is trained on a lot of text data:

- BooksCorpus (800M words)
- English Wikipedia (2.5B words)

Yay, for transfer learning!

BERT-Base model has 12 transformer blocks, 12 attention heads, 110M parameters!

BERT-Large model has 24 transformer blocks, 16 attention heads, 340M parameters!

Takeaway: BERT is incredible for learning contextaware representations of words and using transfer learning for other tasks (e.g., classification).

Can't generate new sentences though, due to no decoders.

ran

What if we want to generate a new output sequence?

GPT-2 model to the rescue!

Generative Pre-trained Transformer 2

GPT-2 (a Transformer variant)

- GPT-2 uses only Transformer Decoders (no Encoders) to generate new sequences
- As it processes each word/token, it cleverly masks the "future" words and conditions itself on the previous words
- Can generate text from scratch or from a starting sequence.
- Easy to fine-tune on your own dataset (language)



GPT-2 (a Transformer variant)

GPT-2 is:

- trained on 40GB of text data (8M webpages)!
- **1.5B parameters**

GPT-3 is an even bigger version (175B parameters) of GPT-2, but isn't open-source

Yay, for transfer learning!

QUESTIONS?