Deep Learning Lab 12-1: Seq2Seq Learning & Neural Machine Translation

DataLab

Department of Computer Science, National Tsing Hua University, Taiwan

Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Encoder-Decoder Model

- As we have seen in lab11-2 (AdaIN),
 - Encoder: encode the inputs into hidden state/context/code
 - Decoder: the hidden state is passed into the decoder to generate the desired output



Encoder-Decoder Model

- As we have seen in lab11-2 (AdaIN),
 - Encoder: encode the inputs into hidden state/context/code
 - Decoder: the hidden state is passed into the decoder to generate the desired output



Encoder-Decoder Model

- As we have seen in lab11-2 (AdaIN),
 - Encoder: encode the inputs into hidden state/context/code
 - Decoder: the hidden state is passed into the decoder to generate the desired output



Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Sequence-to-Sequence (Seq2Seq)

- Sequence-to-Sequence (Seq2Seq) is an architecture based on the encoder-decoder, which transforms an input sequence to the target sequence
 - Both sequences can have arbitrary lengths
 - Have achieved a lot of success in tasks like machine translation, text summarization, and image captioning
 - Google Translate started using such a model in production in late 2016

Google Translate

- Sequence-to-Sequence (Seq2Seq) is an architecture based on the encoder-decoder, which transforms an input sequence to the target sequence
 - The encoder processes each item in the input sequence, it compiles the information it captures into a vector, called the context. After processing the entire input sequence, the encoder sends the context over to the decoder, which begins producing the output sequence item by item



 In neural machine translation, a sequence is a series of words, processed one after another. The output is, likewise, a series of words:



- The context is a vector (an array of numbers, basically) in the case of machine translation. The encoder and decoder tend to both be recurrent neural networks (RNNs)
 - You can set the size of the context vector when you set up your model. It is basically the number of hidden units in the encoder RNN



 By design, a RNN takes two inputs at each time step: an input (in the case of the encoder, one word from the input sentence), and a hidden state

Recurrent Neural Network



- Since the encoder and decoder are both RNNs, each time step one of the RNNs does some processing, it updates its hidden state based on its inputs and previous inputs it has seen
 - Notice the last hidden state is actually the context we pass along to the decoder



Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Why Attention?

- The fixed-length context vector turned out to be a bottleneck for these types of models. It made it challenging for the models to deal with long sentences
 - It is hard for the fixed-length context vector to store all information as the input sequence getting longer and longer
 - It has often forgotten the earlier part of input once it processed the whole input sequence

- Attention allows the model to focus on the relevant parts of the input sequence as needed
 - The decoder can focus on different part of the input sequence at each time step, in order to make a better prediction

Time step: 7



 First, the encoder passes a lot more data to the decoder.
 Instead of passing the last hidden state of the encoding stage, the encoder passes all the hidden states to the decoder:

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



- 2. Second, an attention decoder does an extra step before producing its output. In order to focus on the parts of the input that are relevant to this decoding time step, the decoder does the following:
 - Look at the set of encoder hidden states it received each encoder hidden states is most associated with a certain word in the input sentence
 - 2. Give each hidden states a score (let's ignore how the scoring is done for now)
 - 3. Multiply each hidden states by its softmax score, thus amplifying hidden states with high scores, and drowning out hidden states with low scores

Attention at time step 4

	_	_		_	_	_			_	_	_	_	 		_	 	 	_	_	_	_				_	_	-
1.1																											I
1.11																											J.
1																											I.
1																											J.
1.11																											I.
1																											I.
1																											
1																											I.
1																											I.
1																											I.
1																											
1																											I.
1																											I.
1																											I.
1																											
1																											I.
1																											I.
i																											I.
1																											
î.																											I.
i i																											I.
Ĩ																											I.
S		 	- 1	 		-	-		 			-	 	-	-	 						-	-	 			J.

- Let us now bring the whole thing together in the following visualization and look at how the attention process works:
 - The attention decoder RNN takes in the embedding of the <END> token, and an initial decoder hidden state
 - 2. The RNN processes its inputs, producing an output and a new hidden state vector (h4). The output is discarded
 - 3. Attention Step: We use the encoder hidden states and the h4 vector to calculate a context vector (C4) for this time step
 - 4. We concatenate h4 and C4 into one vector
 - 5. We pass this vector through a feedforward neural network (one trained jointly with the model)
 - 6. The output of the feedforward neural networks indicates the output word of this time step
 - 7. Repeat for the next time steps

Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



 This is another way to look at which part of the input sentence we're paying attention to at each decoding step:



 It actually learned from the training phase how to align words in that language pair



• Let's dive into math!

• Score:
$$e_{ij} = a(s_{i-1}, h_j)$$

• Score after softmax: $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$

• Context vector:
$$c_i = \sum_{j=1}^{T_{\chi}} \alpha_{ij} h_j$$

• How to calculate the score $e_{ij} = a(s_{i-1}, h_j)$?

• , where $a(\cdot)$ is the dot product in the following example

 $decoder_hidden = [10, 5, 10]$

encoder_hidden score
[0, 1, 1] 15 (= 10×0 + 5×1 + 10×1, the dot product)
[5, 0, 1] 60
[1, 1, 0] 15
[0, 5, 1] 35

• Score after softmax
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$
 is straightforward

$\begin{bmatrix} 0, 1, 1 \end{bmatrix} & 15 & 0 \\ \begin{bmatrix} 5, 0, 1 \end{bmatrix} & 60 & 1 \\ \begin{bmatrix} 1, 1, 0 \end{bmatrix} & 15 & 0 \end{bmatrix}$	encoder_l	hide	den	score	score^
	[0,	1,	1]	15	0
	[5,	0,	1]	60	1
	[1,	1,	0]	15	0

• Finally, we can multiply each hidden state of encoder by its score and sum up the alignment vectors to get the context

vector:
$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

enco	odeı	r s	score	score'	`al	ligr	nment				
[0, [5, [1, [0,	1, 0, 1, 5,	1] 1] 0] 1]	15 60 15 35	0 1 0 0	[0, [5, [0, [0,	0, 0, 0, 0,	0] 1] 0] 0]				
cont	text	t =	[0+5+0)+0, 0+	+0+0-	+0,	0+1+0+0]	=	[5,	0,	1]

• There are many well-known score functions as follows

• h represents hidden state of encoder

• *S* represents decoder hidden states

Name	Alignment score function	Citation		
Content-base attention	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) = ext{cosine}[oldsymbol{s}_t,oldsymbol{h}_i]$	Graves2014		
Additive(*)	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) = oldsymbol{v}_a^ op ext{tanh}(oldsymbol{W}_a[oldsymbol{s}_t;oldsymbol{h}_i])$	Bahdanau2015		
Location- Base	$lpha_{t,i} = ext{softmax}(\mathbf{W}_a s_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	Luong2015		
General	$\operatorname{score}(s_t, h_i) = s_t^\top \mathbf{W}_a h_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer.	Luong2015		
Dot-Product	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) = oldsymbol{s}_t^ opoldsymbol{h}_i$	Luong2015		
Scaled Dot-	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) = rac{oldsymbol{s}_t^{ op}oldsymbol{h}_i}{\sqrt{n}}$	Vaswani2017		
Product(^)	Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.			

- Attention mechanism allows the decoder to focus on various part of input sequence, instead of forcing it to encode all information into one fixed-length vector
 - Pros: model interpretability, better performance
 - Cons: computationally expensive

Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Teacher Forcing

- Teacher forcing is a method for quickly and efficiently training recurrent neural network models that use the ground truth from a prior time step as input
 - Accelerate the convergence speed
 - Stabilize the training process



Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Assignment

- There are two parts in the notebook
 - Part I: neural machine translation
 - Part II: sentiment analysis (assignment)

Overview

- Task: translate the sentence from Chinese to English
- Dataset size: 20289
- Encoder: RNN with GRU cell
- Decoder: RNN with GRU cell
- Attention machanism: Bahdanau Attention

$$score(s_t, h_i) = v_a^T tanh(W_a[s_t; h_i])$$

• It is worth noticing that in GRU, the hidden state and the output are same

 It's a toy model with toy dataset. Here we focus on the main idea behind the model



 It's a toy model with toy dataset. Here we focus on the main idea behind the model



 It's a toy model with toy dataset. Here we focus on the main idea behind the model



 In the plotting function, you need to change the path of the Chinese font. Otherwise, the Chinese character will not be displayed in the plot

```
def plot_attention(attention, sentence, predicted_sentence):
 # you need to change the fname based on your system, and the Chinese can be displayed
 in the plot
 font = FontProperties fname=r"./data/TaipeiSansTCBeta-Regular.ttf", size=14)
```

Sentiment Analysis

- Overview
 - Task: predict whether the comment is positive or negative
 - Dataset: IMDB
 - Dataset size: 50000
 - Encoder: RNN with GRU cell
 - Decoder: 4 fully-connected layers
 - Attention machanism: Luong Attention

	review	sentiment
0	One of the other reviewers has mentioned that	positive
1	A wonderful little production. The	positive
2	I thought this was a wonderful way to spend ti	positive
3	Basically there's a family where a little boy	negative
4	Petter Mattei's "Love in the Time of Money" is	positive

TODO

Implement the Luong Attention, where the formula of the score function is:

$$score(s_t, h_i) = s_t^T W_a h_i$$

- h_i : hidden state of the encoder
- S_t : hidden state of the decoder
- W_a : the trainable weights

Demo

- This simple model achieves ~84.5% accuracy with only 10 epochs! Not bad at all!
- Besides the nice accuracy, let's try to do some more fascinating things. How about visualizing our results?

Demo

Positive

y_true: 1

y_predict: 1

jane austen would definitely of this one paltrow does an awesome job capturing the attitude of emma she is funny without being silly yet elegant she puts on very convincing british accent n ot being british myself maybe m not the best judge but she fooled me she was also excellent in doors sometimes forget she american also brilliant are jeremy northam and sophie thompson and law emma thompson sister and mother as the bates women they nearly steal the show and ms law d oesn even have any lines highly recommended

Negative

y_true: 0

y_predict: 0

reaches the point where they become obnoxious and simply frustrating touch football puzzle fami ly and talent shows are not how actual people behave it almost sickening another big flaw is th e woman carell is supposed to be falling for her in her first scene with steve carell is like w atching stroke victim trying to be what imagine is supposed to be unique and original in this w oman comes off as mildly retarded it makes me think that this movie is taking place on another planet left the theater wondering what just saw after thinking further don think it was much

Requirement

- The accuracy should be at least **0.80**
- Show the **10-most-focused words** in the sentence
- Only need to show the **first 10 results** in the test data
- Submit on eeclass your code file Lab12-1_{student id}.ipynb
- No need to submit the checkpoints file, but you should show the results in the notebook
- Deadline: 2023-11-23 (Thu) 23:59

Outline

- Encoder-Decoder Model
- Sequence-to-Sequence (Seq2Seq)
- Attention Mechanism
- Teacher Forcing
- Assignment
- Reference

Reference

- <u>Sequence to Sequence Learning with Neural Networks, I. Sutskever et</u> <u>al., NeurIPS'14</u>
- Learning Phrase Representations using RNN Encoder—Decoder for Statistical Machine Translation, K. Cho et al., EMNLP'14
- <u>Neural Machine Translation by Jointly Learning to Align and Translate,</u> <u>D. Bahdanau et al., ICLR'15</u>
- <u>Effective Approaches to Attention-based Neural Machine Translation,</u> <u>M. T. Luong, EMNLP'15</u>
- Attention Is All You Need, Google Brain, NeurIPS' 17
- <u>Visualizing A Neural Machine Translation Model</u>
- <u>The Illustrated Transformer</u>