

LI7: Neural Machine Translation - 2

Spring 2021

COS 484/584

Last time: Sequence to Sequence learning (Seq2seq)



- Encode entire input sequence into a single vector (using an RNN)
- Decode one word at a time (again, using an RNN!)

(Sutskever et al., 2014)



How seq2seq changed the MT landscape

seq2seq Search term			
United States 💌	2004 - present 💌	All categories 🔻	
Interest over time			
100			
100			
75			
50			
25			
			Not
Jan 1, 2004		Oct 1, 2008	

+ Compare

Web Search 🔻







(source: Rico Sennrich)

RESEARCH > PUBLICATIONS >

Google's Neural Machine **Translation System: Bridging** the Gap between Human and **Machine Translation**

Table 10: Mean of side-by-side scores on production data							
PBMT	GNMT	Human	Relative				
			Improvement				
4.885	5.428	5.504	87%				
4.932	5.295	5.496	64%				
4.035	4.594	4.987	58%				
4.872	5.187	5.372	63%				
5.046	5.343	5.404	83%				
3.694	4.263	4.636	60%				
	of side-by- PBMT 4.885 4.932 4.035 4.872 5.046 3.694	of side-by-side score PBMT GNMT 4.885 5.428 4.932 5.295 4.035 4.594 4.872 5.187 5.046 5.343 3.694 4.263	$\begin{array}{c c c c c c c c c c c c c c c c c c c $				

Table 10: Mean of side-by-side scores on production data							
	PBMT	GNMT	Human	Relative			
				Improvement			
$\mathbf{English} \to \mathbf{Spanish}$	4.885	5.428	5.504	87%			
$\mathbf{English} \to \mathbf{French}$	4.932	5.295	5.496	64%			
$\mathbf{English} \to \mathbf{Chinese}$	4.035	4.594	4.987	58%			
$\text{Spanish} \to \text{English}$	4.872	5.187	5.372	63%			
$French \rightarrow English$	5.046	5.343	5.404	83%			
$\mathbf{Chinese} \to \mathbf{English}$	3.694	4.263	4.636	60%			

(Wu et al., 2016)



Versatile seq2seq

- Seq2seq finds applications in many other tasks!
- Any task where inputs and outputs are sequences of words/ characters
 - Summarization (input text \rightarrow summary)
 - Dialogue (previous utterance \rightarrow reply)
 - > Parsing (sentence \rightarrow parse tree in sequence form)
 - Question answering (context+question \rightarrow answer)

Issues with vanilla seq2seq



- A single encoding vector, *h^{enc}*, needs to capture all the information about source sentence
- Longer sequences can lead to vanishing gradients
- Model may "overfit" to training sequences

Issues with vanilla seq2seq



- Longer sequences can lead to vanishing gradients
- Model may "overfit" to training sequences

> A single encoding vector, h^{enc}, needs to capture all the information about source sentence



$$\mathbf{a} = (3, 4, 2, 1)^\top$$

Remember alignments?



 $\mathbf{a} = (1, 2, 3, 0, 4)^{\top}$

Attention

- The neural MT equivalent of alignment models
- Key idea: At each time step during decoding, focus on a particular part of source sentence
 - This depends on the decoder's current hidden state h^{dec} (i.e. an idea of what you are trying to decode)
 - Usually implemented as a probability distribution over the hidden states of the encoder (h_i^{enc})

Seq2seq with attention



Decoder RNN

(slide credit: Abigail See)



Take softmax to turn the scores into a probability distribution





Use the attention distribution to take a weighted sum of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Decoder RNN







Computing attention



• Encoder hidden states: $h_1^{enc}, \ldots, h_n^{enc}$

• Decoder hidden state at time t: h_t^{dec}

First, get attention scores for this time step of decoder (we'll define g soon): $e^{t} = [g(h_1^{enc}, h_t^{dec}), \dots, g(h_n^{enc}, h_t^{dec})]$

Obtain the attention distribution using softmax:

 $\alpha^{t} = \text{softmax} (e^{t}) \in \mathbb{R}^{n}$

Compute weighted sum of encoder hidden states:

$$a_t = \sum_{i=1}^n \alpha_i^t h_i^{enc} \in \mathbb{R}^h$$

Finally, concatenate with decoder state and pass on to output layer:

	Encoder hidden state	
Je	hidden state #1	
suis	hidden state #2	
étudiant	hidden state #3	

(credits: Jay Alammar)



Types of attention

- **Dot-product attention** (assumes equal dimensions for h^{enc} and h^{dec}): 1. $e_i = g(h_i^{enc}, h^{dec}) = (h^{dec})^T h_i^{enc} \in \mathbb{R}$
- 2. Multiplicative attention:
- 3. Additive attention:

 $g(h_i^{enc}, h^{dec}) = v^T \tanh(W_1 h_i^{enc} + W_2 h^{dec}) \in \mathbb{R}$ where W_1, W_2 are weight matrices (learned) and v is a weight vector (learned)

Assume encoder hidden states $h_1^{enc}, h_2^{enc}, \ldots, h_n^{enc}$ and a decoder hidden state h^{dec}

 $g(h_i^{enc}, h^{dec}) = (h^{dec})^T W h_i^{enc} \in \mathbb{R}$, where W is a weight matrix (learned)



A) the

B) cat

C) sat

Assuming we use dot product attention, which input word will have the highest attention value at current time step?

Dot-product attention:

 $g(h_i^{enc}, h^{dec}) = h^{dec} \cdot h^{enc}$

the -> -0.05 + 0.02cat -> -0.02 + 0.08sat -> 0.01 + 0.04







- A) the
- B) cat
- C) sat

Multiplicative attention:

 $g(h_i^{enc}, h^{dec}) = (h^{dec})^T W h_i^{enc}$

What if we use multiplicative attention with $W = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$? Which input word will have the highest attention value at current time step?

> the -> -0.05 cat -> -0.02 sat -> 0.01



Which value of W in multiplicative attention will provide the same word with highest attention value as dot-product attention?

Multiplicative

attention:

 $g(h_i^{enc}, h^{dec}) = (h^{dec})^T W h_i^{enc}$

$$A) W = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

B)
$$W = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$
 C) both



Attention improves translation

System

Winning WMT'14 system – phrase-based + Existing NMT systems

RNNsearch (Jean et al., 2015)

RNNsearch + unk replace (Jean et al., 2015)

RNNsearch + unk replace + large vocab + en

Our NMT systems

Base

Base + reverse

Base + reverse + dropout

- Base + reverse + dropout + global attention (
- Base + reverse + dropout + global attention (
- Base + reverse + dropout + local-p attention
- Base + reverse + dropout + local-p attention

Ensemble 8 models + unk replace

	Ppl	BLEU
large LM (Buck et al., 2014)		20.7
		16.5
		19.0
semble 8 models (Jean et al., 2015)		21.6
	10.6	11.3
	9.9	12.6 (+1.3)
	8.1	14.0 (+ <i>1.4</i>)
location)	7.3	16.8 (+2.8)
location) + feed input	6.4	18.1 (+1.3)
(general) + feed input	5.0	19.0 (+0.9)
(general) + feed input + unk replace	5.9	20.9 (+1.9)
		23.0 (+2.1)

(Luong et al., 2015)

Visualizing attention



(credits: Jay Alammar)





Going all in on attention

- More recent models (e.g. Transformer, Vaswani et al., 2017) have replaced RNNs entirely with attention mechanisms
- Theoretically limiting (since recurrence can help handle arbitrarily long sequences)
- Huge gains in practical performance



WMT 2014, English-German

Issues with vanilla seq2seq



- A single encoding vector, *h^{enc}*, needs to capture all the information about source sentence
- Longer sequences can lead to vanishing gradients
- Model may "overfit" to training sequences

Dropout

(a) Standard Neural Net

 \otimes \otimes \otimes \otimes $\langle X \rangle$

(b) After applying dropout.

- Form of regularization for RNNs (and any NN in general)
- Idea: "Handicap" NN by removing hidden units stochastically
 - > set each hidden unit in a layer to 0 with probability p during training (p = 0.5 usually works well)
- - ▶ scale outputs by 1/(1-p)
 - hidden units forced to learn more general patterns and improve redundancy
- Test time: Simply compute identity

(Srivastava et al., 2014)







		BLEU					
ID	System	5k	10k	20k	40k	80k	165k
1	Transformer-big	3.3	3.4	4.3	4.7	5.1	5.5
2	Transformer-base	8.3	11.9	16.8	23.2	28.0	32.1
3	2 + feed-forward dimension (2048 \rightarrow 512)	8.8	12.0	16.7	22.3	27.7	31.7
4	3 + attention heads $(8 \rightarrow 2)$	9.2	12.7	19.0	23.6	28.7	32.3
5	4 + dropout (0.1 \rightarrow 0.3)	10.6	17.0	21.9	26.7	31.0	33.4
6	$5 + \text{layers} (6 \rightarrow 5)$	10.9	16.9	21.9	26.0	30.2	33.0
7	6 + label smoothing $(0.1 \rightarrow 0.6)$	11.3	16.5	22.0	26.9	30.4	33.3
8	7 + decoder layerDrop (0 \rightarrow 0.3)	12.9	17.3	22.5	26.9	30.3	33.1
9	8 + target word dropout (0 \rightarrow 0.1)	13.7	18.1	23.1	27.0	30.7	33.0
10	9 + activation dropout (0 \rightarrow 0.3)	14.3	18.3	23.6	27.4	30.4	32.6

Table 2: Results of Transformer optimized on the 5k dataset for different subsets and full corpus of IWSLT14 German \rightarrow English. Averages over three runs from three different samples are reported.

(Araabi and Monz, 2020)



Other challenges with NMT

- Out-of-vocabulary words
- Low-resource languages
- Long-term context
- Common sense knowledge (e.g. hot dog, paper jam)
- Fairness and bias
- Uninterpretable

Massively multilingual MT



- Massive improvements on low-resource languages

(Arivazhagan et al., 2019)

Train a single neural network on 103 languages paired with English (remember Interlingua?)

Bilingual Baselines →





ENGLISH HINDI - DETECTED ←



vo sundar hai. vo buddhimaan hai. vo padhaakoo hai. vo vyast hai. vo ameer hai.

Did you mean: वो सुंदर है. वो बुद्धिमान है. वो पढ़ाई
है. वो व्यस्त है. वो अमीर है.



70 / 5000

She is beautiful He is intelligent 🕸 He is a nerd. He is busy He is rich





Bias and Fairness



- NMT systems suffer from issues of systematic bias (e.g. gender)
- Evident when translating from/to a language with gender-specific (or genderagnostic) terms
- Models learn (and amplify) stereotypes from data





Measuring bias in MT



- non-stereotypical gender roles
 - e.g. "The doctor asked the nurse to help her in the operation"
- Systems consistently performed worse on non-stereotypical gender translation

WinoMT: Stanovsky et al. (2019) use coreference resolution to construct a dataset of

Source	[Target lang.] Predicted translation	Phenomenon		
The janitor does not like the baker because she always messes up the kitchen.	[ES] Al conserje no le gusta <mark>el panadero</mark> porque ella siempre desordena la cocina.	Biased translation, giving "baker" a male inflection, with a mismatched pro- noun reference.		
The janitor does not like the pretty baker because she always messes up the kitchen.	[ES] Al conserje no le gusta la panadera bonita porque ella siempre desordena la cocina.	Adding a stereotypically female adjec- tive "fixes" the translation.		
The counselor asked the guard a few questions and praised her for the good work.	[FR] Le conseiller a posé quelques ques- tions à la garde et l'a louée pour le bon travail.	French uses "garde" for both male and female guards, allowing for a more direct translation from English.		

Table 5: Examples of Google Translate's output for different sentences in the WinoMT corpus. Words in blue, red, and orange indicate male, female and neutral entities, respectively.

(Stanovsky et al. 2019)

Mitigating bias

- systems
- w/ or w/o stereotypes)

	WMT Data Systems					
uit for her interview.			Acc.	$\Delta \mathbf{G}$	$\Delta \mathbf{S}$	M:H
	(*)	Baseline	66.7	10.2	14.4	2.6
	Ð	TGA Oracle	89.0	-4.7	1.7	1
uit for ner interview.	Ľ.	TGA HuggingFace	77.6	-0.1	11.9	1.6
UUFU	щ	TGA AllenNLP	81.5	-2.0	11.1	1.4
	EN-FR	Baseline	48.6	29.8	11.8	5.5
ams jauns uzvalks		TGA Oracle	81.5	1.4	2.8	1.2
anis jauns uzvaiks.		TGA HuggingFace	67.8	4.9	12.4	2
		TGA AllenNLP	74.4	1.6	10.1	1.6

Stafanovics et al. (2020) use word-level annotations of subject's gender to train NMT

• TGA (target gender annotations) help reduce gender bias (∇G = diff. in F1 between sentences with male and female antecedents, $\nabla S = \text{diff.}$ in accuracy between sentences

:F .4 .5 .2 2

Anonymous feedback form: https://forms.gle/7BxYDUTebogndJQE8