

COS 484: Natural Language Processing

Neural Machine Translation

Fall 2019

Announcements

- Sign up for project meetings ASAP
 - Link will be posted on Piazza by 3pm today
 - Will be held this week

Last time



- Statistical MT
 - Word-based
 - Phrase-based
 - Syntactic



Neural Machine Translation

- A single neural network is used to translate from source to target
- Architecture: Encoder-Decoder
 - Two main components:
 - Encoder: Convert source sentence (input) into a vector/matrix
 - Decoder: Convert encoding into a sentence in target language (output)

Recall: RNNs

 $\mathbf{h}_t = g(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b}) \in \mathbb{R}^d$



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!)
- Beam search for better inference
- Learning is not trivial! (vanishing/exploding gradients)

(Sutskever et al., 2014)

Sentence: This cat is cute



Sentence: This cat is cute



Sentence: This cat is cute













• A conditioned language model



Seq2seq training

- Similar to training a language model!
- Minimize cross-entropy loss:

$$\sum_{t=1}^{T} -\log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- Back-propagate gradients through both decoder and encoder
- Need a really big corpus

36M sentence pairs Russian: Машинный перевод - это круто! figlish: Machine translation is cool!

Seq2seq training



⁽slide credit: Abigail See)

Greedy decoding



- Compute argmax at every step of decoder to generate word
- What's wrong?

Exhaustive search?

Find arg max
$$P(y_1, \dots, y_T | x_1, \dots, x_n)$$

 y_1, \dots, y_T

- Requires computing all possible sequences
 - $O(V^T)$ complexity!
 - Too expensive

A middle ground: Beam search

- Key idea: At every step, keep track of the k most probable partial translations (hypotheses)
- Score of each hypothesis = log probability

$$\sum_{t=1}^{j} \log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- Not guaranteed to be optimal
- More efficient than exhaustive search

Beam size = k = 2. Blue numbers = $score(y_1, \ldots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i|y_1, \ldots, y_{i-1}, x)$



Beam size = k = 2. Blue numbers = $score(y_1, \ldots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i|y_1, \ldots, y_{i-1}, x)$





Backtrack



- Different hypotheses may produce $\langle e \rangle$ (end) token at different time steps
 - When a hypothesis produces $\langle e \rangle$, stop expanding it and place it aside
- Continue beam search until:
 - All k hypotheses produce $\langle e \rangle$ OR
 - Hit max decoding limit T
- Select top hypotheses using the normalized likelihood score

$$\frac{1}{T} \sum_{t=1}^{T} \log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

Otherwise shorter hypotheses have higher scores

NMT vs SMT



Cons

NMT vs SMT

Pros

- Better performance
 - Fluency
 - Longer context
- Single NN optimized end-toend
- Less engineering
- Works out of the box for many language pairs

Cons

- Requires more data and compute
- Less interpretable
 - Hard to debug
- Uncontrollable
 - Heavily dependent on data could lead to unwanted biases
- More parameters

How seq2seq changed the MT landscape



MT Progress



(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 score	s on prou	uction data
	PBMT	GNMT	Human	Relative
				Improvement
$English \rightarrow Spanish$	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%

Table 10. Moon of side by side secres on production date

(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
- Overfitting

Remember alignments?





Τ

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

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 (i.e. notion of what you are trying to decode)
 - Usually implemented as a probability distribution over the hidden states of the encoder (h_i^{enc})











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 (we will see what g is soon!): $e^t = [g(h_1^{enc}, h_t^{dec}), \dots, g(h_n^{enc}, h_t^{dec})]$
- Obtain the attention distribution using softmax:

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

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: $[a_t; h_t^{dec}] \in \mathbb{R}^{2h}$

Types of attention

- Assume encoder hidden states h_1, h_2, \ldots, h_n and decoder hidden state z
- 1. Dot-product attention (assumes equal dimensions for *a* and *b*: $e_i = g(h_i, z) = z^T h_i \in \mathbb{R}$
- 2. Multiplicative attention:

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

3. Additive attention:

 $g(h_i, z) = v^T \tanh(W_1 h_i + W_2 z) \in \mathbb{R}$

where W_1 , W_2 are weight matrices and v is a weight vector

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

Overfitting

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
- Test time: Simply compute identity



(a) Standard Neural Net



(b) After applying dropout.

Existing 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

02 -	Slavic	u be	Turkic • ^{kk}	bn gur Indo	o-Aryan ian
0.01 -	• ^{bg}	uk Baltic	hy ky ky ja _{fa}	si ps am	d
0-	cs. ^{sk} pl hr sl	ethu es ^{phu} Uralic	iw tr ar az th yi	km my Kra-Dai	
-0.01 -	sv fr Plda gl	de ro sq eo af	is ga vi mt ht b	Niger-	Congo ^{xh}
-	Romance	Germanic	Austronesian	su ceb	st sh ha mi ig ny haw

- Train a single neural network on 103 languages paired with English (remember Interlingua?)
- Massive improvements on low-resource languages

(Arivazhagan et al., 2019)