Lecture 13: Machine Translation II

(many slides from Greg Durrett)

Alan Ritter

Neural MT Details

Sutskever seq2seq paper: first major application of LSTMs to NLP

Sutskever et al. (2014)



- Sutskever seq2seq paper: first major application of LSTMs to NLP
- Basic encoder-decoder with beam search





- Sutskever seq2seq paper: first major application of LSTMs to NLP
- Basic encoder-decoder with beam search



	test BLEU score (ntst14)
	28.45
	33.30
12	26.17
12	30.59
	•
size 12	34.81

Sutskever et al. (2014)



- Sutskever seq2seq paper: first major application of LSTMs to NLP
- Basic encoder-decoder with beam search



test BLEU score (ntst14)
28.45
33.30
26.17
30.59
34.81

Sutskever et al. (2014)



- and copying for rare words



Better model from seq2seq lectures: encoder-decoder with attention

distribution over vocab + copying

12M sentence pairs

12M sentence pairs

Classic phrase-based system: ~33 BLEU, uses additional target-language data



12M sentence pairs

Classic phrase-based system: ~33 BLEU, uses additional target-language data Rerank with LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014)



- 12M sentence pairs
- Classic phrase-based system: ~33 BLEU, uses additional target-language data Rerank with LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014) Sutskever+ (2014) seq2seq single: **30.6** BLEU



- 12M sentence pairs
- Classic phrase-based system: ~33 BLEU, uses additional target-language data
 - Rerank with LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014)
- Sutskever+ (2014) seq2seq single: **30.6** BLEU
- Sutskever+ (2014) seq2seq ensemble: **34.8** BLEU



- 12M sentence pairs
- Classic phrase-based system: ~33 BLEU, uses additional target-language data
 - Rerank with LSTMs: **36.5** BLEU (long line of work here; Devlin+ 2014)
- Sutskever+ (2014) seq2seq single: **30.6** BLEU
- Sutskever+ (2014) seq2seq ensemble: 34.8 BLEU
- Luong+ (2015) seq2seq ensemble with attention and rare word handling: **37.5** BLEU
- But English-French is a really easy language pair and there's tons of data for it! Does this approach work for anything harder?





Results: WMT English-German

- 4.5M sentence pairs
- Classic phrase-based system: **20.7** BLEU
- Luong+ (2014) seq2seq: **14** BLEU
- Luong+ (2015) seq2seq ensemble with rare word handling: **23.0** BLEU
- languages

Not nearly as good in absolute BLEU, but not really comparable across

Results: WMT English-German

- 4.5M sentence pairs
- Classic phrase-based system: **20.7** BLEU
- Luong+ (2014) seq2seq: **14** BLEU
- Luong+ (2015) seq2seq ensemble with rare word handling: 23.0 BLEU
- languages
- French, Spanish = easiest German, Czech = harder

Not nearly as good in absolute BLEU, but not really comparable across

Japanese, Russian = hard (grammatically different, lots of morphology...)



MT Examples

src	In einem Interview sagte Bloom jedoch
ref	However, in an interview, Bloom has s
best	In an interview, however, Bloom said t
base	However, in an interview, Bloom said

- best = with attention, base = no attention
- phrase-based doesn't do this

, dass er und Kerr sich noch immer lieben .

said that he and *Kerr* still love each other.

that he and *Kerr* still love .

that he and **Tina** were still $\langle unk \rangle$.

NMT systems can hallucinate words, especially when not using attention

Luong et al. (2015)



MT Examples

src	Wegen der von Berlin und der Europäis
	Verbindung mit der Zwangsjacke, in die
	ten an der gemeinsamen Währung genötig
	Europa sei zu weit gegangen
ref	The austerity imposed by Berlin and the
	imposed on national economies through a
	to think Project Europe has gone too far.
best	Because of the strict austerity measures
	connection with the straitjacket in which
	the common currency, many people belie
base	Because of the pressure imposed by the E
	with the strict austerity imposed on the
	many people believe that the European pro-

best = with attention, base = no attention

schen Zentralbank verhängten strengen Sparpolitik in e die jeweilige nationale Wirtschaft durch das Festhalgt wird, sind viele Menschen der Ansicht, das Projekt

European Central Bank, coupled with the straitjacket dherence to the common currency, has led many people

imposed by Berlin and the European Central Bank in the respective national economy is forced to adhere to eve that the European project has gone too far. uropean Central Bank and the Federal Central Bank e national economy in the face of the single currency, oject has gone too far.

Luong et al. (2015)





MT Examples

~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~					
Source	such changes in reaction conditions include, but are not limited to,				
	an increase in temperature or change in ph .				
Reference	所(such) 述(said) 反 应(reaction) 条 件(condition) 的(of)				
	改 变(change) 包 括(include) 但(but) 不(not) 限 于(limit)				
	温度(temperature) 的(of) 增加(increase) 或(or) pH 值(value) 的(of) 改变(change) 。				
PBMT	中(in) 的(of) 这种(such) 变化(change) 的(of) 反应(reaction) 条				
	件(condition) 包括(include) , 但(but) 不(not) 限于(limit) ,				
增加(increase)的(of)温度(temperature)或(or)pH变化(change)。					
NMT	这种(such) 反应(reaction) 条件(condition) 的(of) 变化(change) 包括(include) 但(but) 不(not)				
	限于(limit) pH 或(or) pH 的(of) 变化(change)。				

- NMT can repeat itself if it gets confused (pH or pH)
- Phrase-based MT often gets chunks right, may have more subtle ungrammaticalities

Zhang et al. (2017)



Use Huffman encoding on a corpus, keep most common k (~10,000) character sequences for source and target





Use Huffman encoding on a corpus, keep most common k (~10,000) character sequences for source and target



Captures common words and parts of rare words



Use Huffman encoding on a corpus, keep most common k (~10,000) character sequences for source and target



- Captures common words and parts of rare words
- Subword structure may make it easier to translate



Use Huffman encoding on a corpus, keep most common k (~10,000) character sequences for source and target



- Captures common words and parts of rare words
- Subword structure may make it easier to translate
- Model balances translating and transliterating without explicit switching Wu et al. (2016)



- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters



- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters
- for i in range(num_merges): pairs = get_stats(vocab) best = max(pairs, key=pairs.get) vocab = merge_vocab(best, vocab)



- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters
- for i in range(num_merges): pairs = get_stats(vocab) best = max(pairs, key=pairs.get) vocab = merge_vocab(best, vocab)
 - Count bigram character cooccurrences





- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters
- for i in range(num_merges): Count bigram character cooccurrences pairs = get_stats(vocab) best = max(pairs, key=pairs.get) Merge the most frequent pair of vocab = merge_vocab(best, vocab)
 - adjacent characters





- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters
- for i in range(num_merges): Count bigram character cooccurrences pairs = get_stats(vocab) best = max(pairs, key=pairs.get) Merge the most frequent pair of vocab = merge_vocab(best, vocab)

- Final size = initial vocab + num merges. Often do 10k 30k merges
- adjacent characters





- Simpler procedure, based only on the dictionary
- Input: a dictionary of words represented as characters
- for i in range(num_merges): Count bigram character cooccurrences pairs = get_stats(vocab) best = max(pairs, key=pairs.get) Merge the most frequent pair of vocab = merge_vocab(best, vocab)

- Final size = initial vocab + num merges. Often do 10k 30k merges
- Most SOTA NMT systems use this on both source + target

adjacent characters







Google's NMT System



Google's NMT System



English-French:

- Google's phrase-based system: 37.0 BLEU Luong+ (2015) seq2seq ensemble with rare word handling: 37.5 BLEU Google's 32k word pieces: 38.95 BLEU

Google's NMT System



English-French:

Google's phrase-based system: 37.0 BLEU Luong+ (2015) seq2seq ensemble with rare word handling: 37.5 BLEU Google's 32k word pieces: 38.95 BLEU

English-German:

Google's phrase-based system: 20.7 BLEU Luong+ (2015) seq2seq ensemble with rare word handling: 23.0 BLEU Google's 32k word pieces: 24.2 BLEU

Google's NMT System



Human Evaluation (En-Es)

200

100

0

Similar to human-level 400 performance on English-Spanish 300 Count (total 500)



PBMT - GNMT - Human



Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard j
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard
	coincée dans la carrière

Google's NMT System

dog walker trapped in the quarry par un promeneur de chien piégé dans la carrière l par un traîneau à chiens piégé dans la carrière.

l par une personne qui promenait son chien





Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard
	coincée dans la carrière

Google's NMT System

dog walker trapped in the quarry par un promeneur de chien piégé dans la carrière l par un traîneau à chiens piégé dans la carrière.

l par une personne qui promenait son chien





Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard
	coincée dans la carrière

Google's NMT System





Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard
	coincée dans la carrière

Google's NMT System





Source	She was spotted three days later by a
PBMT	Elle a été repéré trois jours plus tard
GNMT	Elle a été repérée trois jours plus tard
Human	Elle a été repérée trois jours plus tard
	coincée dans la carrière

Google's NMT System





do the same?

Classical MT methods used a bilingual corpus of sentences B = (S, T) and a large monolingual corpus T' to train a language model. Can neural MT



- do the same?
- Approach 1: force the system to generate T' as targets from null inputs

Classical MT methods used a bilingual corpus of sentences B = (S, T) and a large monolingual corpus T' to train a language model. Can neural MT



- do the same?
- Approach 1: force the system to generate T' as targets from null inputs

Classical MT methods used a bilingual corpus of sentences B = (S, T) and a large monolingual corpus T' to train a language model. Can neural MT



- do the same?
- Approach 1: force the system to generate T' as targets from null inputs

Classical MT methods used a bilingual corpus of sentences B = (S, T) and a large monolingual corpus T' to train a language model. Can neural MT

> Approach 2: generate synthetic sources with a T->S machine translation system (backtranslation)





- do the same?
- Approach 1: force the system to generate T' as targets from null inputs

Classical MT methods used a bilingual corpus of sentences B = (S, T) and a large monolingual corpus T' to train a language model. Can neural MT

> Approach 2: generate synthetic sources with a T->S machine translation system (backtranslation)







name	training		BLEU			
	data	instances	tst2011	tst2012	tst2013	tst2014
baseline (Gülçehre et al., 2015)			18.4	18.8	19.9	18.7
deep fusion (Gülçehre et al., 2015)		20.2	20.2	21.3	20.6	
baseline	parallel	7.2m	18.6	18.2	18.4	18.3
parallel _{synth}	parallel/parallel _{synth}	6m/6m	19.9	20.4	20.1	20.0
Gigaword _{mono}	parallel/Gigaword _{mono}	7.6m/7.6m	18.8	19.6	19.4	18.2
Gigaword _{synth}	parallel/Gigaword _{synth}	8.4m/8.4m	21.2	21.1	21.8	20.4

- Gigaword: large monolingual English corpus
- parallel_{synth}: backtranslate training data; makes additional noisy source sentences which could be useful



Transformers for MT

Each word forms a "query" which then computes attention over each word

the movie was great



Each word forms a "query" which then computes attention over each word







Each word forms a "query" which then computes attention over each word







Each word forms a "query" which then computes attention over each word









Each word forms a "query" which then computes attention over each word

 $\alpha_{i,j} = \operatorname{softmax}(x_i^{\top} x_j)$ scalar









Each word forms a "query" which then computes attention over each word

$$lpha_{i,j} = \operatorname{softmax}(x_i^{ op} x_j)$$
 scalar $x_i' = \sum_{j=1}^n lpha_{i,j} x_j$ vector = sum of scalar



Each word forms a "query" which then computes attention over each word

$$lpha_{i,j} = \operatorname{softmax}(x_i^{ op} x_j)$$
 scalar $x_i' = \sum_{i=1}^n lpha_{i,j} x_j$ vector = sum of scalar

Multiple "heads" analogous to different convolutional filters. Use

parameters W_k and V_k to get different attention values + transform vectors

Each word forms a "query" which then computes attention over each word

$$lpha_{i,j} = ext{softmax}(x_i^ op x_j)$$
 scalar $x_i' = \sum_{j=1}^n lpha_{i,j} x_j$ vector = sum of scalar

Multiple "heads" analogous to different convolutional filters. Use

$$\alpha_{k,i,j} = \operatorname{softmax}(x_i^\top W_k x_j)$$

parameters W_k and V_k to get different attention values + transform vectors

Each word forms a "query" which then computes attention over each word

$$lpha_{i,j} = ext{softmax}(x_i^ op x_j)$$
 scalar $x_i' = \sum_{j=1}^n lpha_{i,j} x_j$ vector = sum of scalar

Multiple "heads" analogous to different convolutional filters. Use

$$\alpha_{k,i,j} = \operatorname{softmax}(x_i^\top W_k x_j) \quad x'_{k,i} = \sum_{j=1}^n \alpha_{k,i,j} V_k x_j$$

Vaswani et al. (20)

parameters W_k and V_k to get different attention values + transform vectors

Each word forms a "query" which then computes attention over each word

$$lpha_{i,j} = ext{softmax}(x_i^ op x_j)$$
 scalar $x_i' = \sum_{j=1}^n lpha_{i,j} x_j$ vector = sum of scalar

Multiple "heads" analogous to different convolutional filters. Use

$$\alpha_{k,i,j} = \operatorname{softmax}(x_i^\top W_k x_j) \quad x'_{k,i} = \sum_{j=1}^n \alpha_{k,i,j} V_k x_j$$

Vaswani et al. (20)

parameters W_k and V_k to get different attention values + transform vectors

Transformers

- a one-hot vector

Augment word embedding with position embeddings, each dim is a sine/cosine wave of a different frequency. Closer points = higher dot products

Works essentially as well as just encoding position as Vaswani et al. (2017)

Encoder and decoder are both transformers

Decoder consumes the previous generated token (and attends to input), but has no recurrent state

Big = 6 layers, 1000 dim for each token, 16 heads, base = 6 layers + other params halved

Transformers

BLEU			
EN-DE	EN-FR		
23.75			
	39.2		
24.6	39.92		
25.16	40.46		
26.03	40.56		
	40.4		
26.30	41.16		
26.36	41.29		
27.3	38.1		
28.4	41.8		

Visualization

lt	<u>s</u>	. <mark>с</mark>	this	spirit	that	a	majority	of	American	governments	have	passed	New	Swe
Ħ	<u>s</u>	. 드	this	spirit	that	a	majority	of	American	governments	have	passed	New	SWE

Visualization

- Can build MT systems with LSTM encoder-decoders, CNNs, or transformers
- Word piece / byte pair models are really effective and easy to use
- State of the art systems are getting pretty good, but lots of challenges remain, especially for low-resource settings
- Next time: pre-trained transformer models (BERT), applied to other tasks