# Natural Language Processing with Deep Learning CS224N/Ling284



#### Lecture 10: Machine Translation, Sequence-to-sequence and Attention

**Abigail See** 

#### **Overview**

Today we will:

• Introduce a <u>new task</u>: Machine Translation

is the primary use-case of

• Introduce a <u>new neural architecture</u>: sequence-to-sequence

is improved by

• Introduce a <u>new neural technique</u>: attention

# **Machine Translation**

<u>Machine Translation (MT)</u> is the task of translating a sentence x from one language (the source language) to a sentence y in another language (the target language).

*x:* L'homme est né libre, et partout il est dans les fers

*y:* Man is born free, but everywhere he is in chains

# **1950s: Early Machine Translation**

Machine Translation research began in the early 1950s.

 Mostly Russian → English (motivated by the Cold War!)



Source: https://youtu.be/K-HfpsHPmvw

- Systems were mostly rule-based, using a bilingual dictionary to map Russian words to their English counterparts
  - A cool by-product: Quicksort!

- <u>Core idea</u>: Learn a probabilistic model from data
- Suppose we're translating French  $\rightarrow$  English.
- We want to find best English sentence y, given French sentence x

 $\mathrm{argmax}_y P(y|x)$ 

 Use Bayes Rule to break this down into two components to be learnt separately:



- <u>Question</u>: How to learn translation model P(x|y) ?
- First, need large amount of parallel data (e.g. pairs of human-translated French/English sentences)



- <u>Question</u>: How to learn translation model P(x|y) ?
- First, need large amount of parallel data (e.g. pairs of human-translated French/English sentences)
- Break it down further: we actually want to consider

where *a* is the alignment, i.e. word-level correspondence between French sentence *x* and English sentence *y* 

# What is alignment?

Alignment is the correspondence between particular words in the translated sentence pair.

• Note: Some words have no counterpart



# **Alignment is complex**

Alignment can be one-to-many (these are "fertile" words)



# **Alignment is complex**

Alignment can be many-to-one





# **Alignment is complex**

Alignment can be many-to-many (phrase-level)





- <u>Question</u>: How to learn translation model P(x|y) ?
- First, need large amount of parallel data (e.g. pairs of human-translated French/English sentences)
- Break it down further: we actually want to consider

where *a* is the alignment, i.e. word-level correspondence between French sentence *x* and English sentence *y* 

- We learn P(x, a|y) as a combination of many factors, including:
  - Probability of particular words aligning
    - Also depends on position in sentence
  - Probability of particular words having particular fertility
- <sup>15</sup> Etc.



- We could enumerate every possible y and calculate the probability? → Too expensive!
- <u>Answer</u>: Use a heuristic search algorithm to gradually build up the the translation, discarding hypotheses that are too low-probability

## **Searching for the best translation**



## **Searching for the best translation**





- SMT is a huge research field
- The best systems are extremely complex
  - Hundreds of important details we haven't mentioned here
  - Systems have many separately-designed subcomponents
  - Lots of feature engineering
    - Need to design features to capture particular language phenomena
  - Require compiling and maintaining extra resources
    - Like tables of equivalent phrases
  - Lots of human effort to maintain
    - Repeated effort for each language pair!

# 2014

(dramatic reenactment)

# 2014

Milesearch

#### (dramatic reenactment)

anslation

# What is Neural Machine Translation?

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*
- The neural network architecture is called sequence-to-sequence (aka seq2seq) and it involves *two* RNNs.

# **Neural Machine Translation (NMT)**



# **Neural Machine Translation (NMT)**

- The sequence-to-sequence model is an example of a Conditional Language Model.
  - Language Model because the decoder is predicting the next word of the target sentence *y*
  - Conditional because its predictions are *also* conditioned on the source sentence x
- NMT directly calculates P(y|x):

 $P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots, P(y_T|y_1, \dots, y_{T-1}, x)$ 

Probability of next target word, given target words so far and source sentence x

- **Question**: How to train a NMT system?
- <u>Answer</u>: Get a big parallel corpus...

# **Training a Neural Machine Translation system**



ecoder RNN

# **Better-than-greedy decoding?**

 We showed how to generate (or "decode") the target sentence by taking argmax on each step of the decoder



- This is greedy decoding (take most probable word on each step)
- Problems?

# **Better-than-greedy decoding?**

- Greedy decoding has no way to undo decisions!
  - les pauvres sont démunis (the poor don't have any money)
  - $\rightarrow$  the \_\_\_\_\_
  - $\rightarrow$  the poor \_\_\_\_\_
  - $\rightarrow$  the poor are \_\_\_\_\_
- Better option: use beam search (a search algorithm) to explore several hypotheses and select the best one

## **Beam search decoding**

- Ideally we want to find *y* that maximizes  $P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots, P(y_T|y_1, \dots, y_{T-1}, x)$
- We could try enumerating all  $y \rightarrow$  too expensive!
  - Complexity  $O(V^T)$  where V is vocab size and T is target sequence length
- <u>Beam search</u>: On each step of decoder, keep track of the k most probable partial translations
  - *k* is the beam size (in practice around 5 to 10)
  - Not guaranteed to find optimal solution
  - But much more efficient!

Beam size = 2

<START>

Beam size = 2



Beam size = 2



Beam size = 2



Beam size = 2 always not are poor don't the have people take <START> person poor but а person









# **Advantages of NMT**

Compared to SMT, NMT has many advantages:

- Better performance
  - More fluent
  - Better use of context
  - Better use of phrase similarities
- A single neural network to be optimized end-to-end
  - No subcomponents to be individually optimized
- Requires much less human engineering effort
  - No feature engineering
  - Same method for all language pairs

# **Disadvantages of NMT?**

Compared to SMT:

- NMT is less interpretable
  - Hard to debug
- NMT is difficult to control
  - For example, can't easily specify rules or guidelines for translation
  - Safety concerns!

# How do we evaluate Machine Translation?

**BLEU** (Bilingual Evaluation Understudy)

- BLEU compares the <u>machine-written translation</u> to one or several <u>human-written translation(s)</u>, and computes a <u>similarity</u> score based on:
  - *n*-gram precision (usually up to 3 or 4-grams)
  - Penalty for too-short system translations
- BLEU is useful but imperfect
  - There are many valid ways to translate a sentence
  - So a good translation can get a poor BLEU score because it has low *n*-gram overlap with the human translation <sup>(3)</sup>

# **MT progress over time**

[Edinburgh En-De WMT newstest2013 Cased BLEU; NMT 2015 from U. Montréal]



Source: http://www.meta-net.eu/events/meta-forum-2016/slides/09\_sennrich.pdf

#### NMT: the biggest success story of NLP Deep Learning

Neural Machine Translation went from a fringe research activity in **2014** to the leading standard method in **2016** 

- **2014**: First seq2seq paper published
- **2016**: Google Translate switches from SMT to NMT
- This is amazing!
  - SMT systems, built by hundreds of engineers over many years, outperformed by NMT systems trained by a handful of engineers in a few months

- Nope!
- Many difficulties remain:
  - Out-of-vocabulary words
  - Domain mismatch between train and test data
  - Maintaining context over longer text
  - Low-resource language pairs

- Nope!
- Using common sense is still hard



Open in Google Translate





Feedback

- Nope!
- NMT picks up biases in training data



Didn't specify gender

Source: https://hackernoon.com/bias-sexist-or-this-is-the-way-it-should-be-ce1f7c8c683c

- Nope!
- Uninterpretable systems do strange things

English	Spanish	Japanese	Detect language	•		English	Spanish	Arabic	•	Translate
が ががが ががが ががが ががが ががが ががが ががが	が ががが ががががが ががががが ががががが ががががが ががががが	が がが がががが ががががか ががががか ががががわ	ヾ ヾが ヾがが			But Peel A pain I feel a My stor Strange Having My bac Strong Strong There v It is pro Strong	is strange f mach e feeling a bad ap d gray but burns but burns was a bac one to bur but burni	feeling opearand s d shape rns, but shed	ce but a also a	a bad shape a burn
ががががががががががががががが										

**Source**: http://languagelog.ldc.upenn.edu/nll/?p=35120#more-35120

# **NMT research continues**

NMT is the **flagship task** for NLP Deep Learning

- NMT research has pioneered many of the recent innovations of NLP Deep Learning
- In **2018**: NMT research continues to thrive
  - Researchers have found *many, many* improvements to the "vanilla" seq2seq NMT system we've presented today
  - But one improvement is so integral that it is the new vanilla...

# ATTENTION

# Sequence-to-sequence: the bottleneck problem



**Problems with this architecture?** 

## **Sequence-to-sequence: the bottleneck problem**



## Attention

- Attention provides a solution to the bottleneck problem.
- <u>Core idea</u>: on each step of the decoder, *focus on a particular part* of the source sequence



 First we will show via diagram (no equations), then we will show with equations











Decoder RNN



Decoder RNN





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

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





Decoder RNN





2/15/18





2/15/18









2/15/18





## **Attention: in equations**

- We have encoder hidden states  $h_1, \ldots, h_N \in \mathbb{R}^h$
- On timestep t, we have decoder hidden state  $s_t \in \mathbb{R}^h$
- We get the attention scores  $e^t$  for this step:

$$oldsymbol{e}^t = [oldsymbol{s}_t^Toldsymbol{h}_1, \dots, oldsymbol{s}_t^Toldsymbol{h}_N] \in \mathbb{R}^N$$

• We take softmax to get the attention distribution  $\alpha^t$  for this step (this is a probability distribution and sums to 1)

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

$$\boldsymbol{a}_t = \sum_{i=1}^{h} \alpha_i^t \boldsymbol{h}_i \in \mathbb{R}^h$$

• Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model

$$[oldsymbol{a}_t;oldsymbol{s}_t]\in\mathbb{R}^{2h}$$
 2/15/18

# **Attention is great**

- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we can see what the decoder was focusing on
  - We get alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself



#### Recap

- We learned the history of Machine Translation (MT)
- Since 2014, Neural MT rapidly replaced intricate Statistical MT

 Sequence-to-sequence is the architecture for NMT (uses 2 RNNs)

- Attention is a way to *focus on particular parts* of the input
  - Improves sequence-to-sequence a lot!





## Sequence-to-sequence is versatile!

- Sequence-to-sequence is useful for *more than just MT*
- Many NLP tasks can be phrased as sequence-to-sequence:
  - Summarization (long text  $\rightarrow$  short text)
  - Dialogue (previous utterances  $\rightarrow$  next utterance)
  - Parsing (input text → output parse as sequence)
  - Code generation (natural language  $\rightarrow$  Python code)

## **Next time**

- More types of attention
- More uses for attention
- More advanced sequence-to-sequence techniques