Introduction to Attention

Danna Gurari University of Colorado Boulder Spring 2025



https://dannagurari.colorado.edu/course/neural-networks-and-deep-learning-spring-2025/

Review

- Last lecture:
 - Motivation: numeric representation of natural language
 - Tokenization: how to convert text into discrete units
 - Neural word embeddings: how to create dense representation
 - Programming tutorial
- Assignments (Canvas):
 - Problem set 3 grades are out
 - Review session will be held at 2pm today on Zoom
 - Email all regrade requests to our TA, Nick Cooper (a comment in Canvas is not sufficient)
 - Lab assignment 2 due in one week
- Questions?

Today's Topics

- Motivation: machine neural translation for long sentences
- Decoder: attention
- Encoder
- Performance evaluation
- Final project: ways to find a partner

Today's Topics

- Motivation: machine neural translation for long sentences
- Decoder: attention
- Encoder
- Performance evaluation
- Final project: ways to find a partner

Task: Machine Translation

DETECT LANGUAGE ENGLISH SPANISH FRENCH N	· ↓ ↓	GERMAN ENGLISH SPANISH V	
He loved to eat	×	Er liebte es zu essen	
↓ ↓	15 / 5,000 📼 🔻		n c _q <

Which type of sequence problem is this: one-to-many, many-to-one, or many-to-many?

Pioneering Neural Network Approach



Image source: https://smerity.com/articles/2016/google_nmt_arch.html seq2seq: Sutskever et al. Sequence to Sequence Learning with Neural Networks. Neurips 2014.

Pioneering Neural Network Approach



Image source: https://smerity.com/articles/2016/google_nmt_arch.html seq2seq: Sutskever et al. Sequence to Sequence Learning with Neural Networks. Neurips 2014.

Analysis of Two Models



What performance trend is observed for inputs (source) and outputs (reference) as the number of words in each sentence grows?

Cho et al. On the Properties of Neural Machine Translation: Encoder–Decoder Approaches. SSST 2014

Analysis of Two Models



Performance drops for longer sentences!

Cho et al. On the Properties of Neural Machine Translation: Encoder–Decoder Approaches. SSST 2014





Image source: https://smerity.com/articles/2016/google_nmt_arch.html seq2seq: Sutskever et al. Sequence to Sequence Learning with Neural Networks. Neurips 2014.



Image source: https://smerity.com/articles/2016/google_nmt_arch.html



Image source: https://smerity.com/articles/2016/google_nmt_arch.html

Decoder decides which inputs are needed for prediction at each time step; e.g., "hard attention" focuses on one input



Note: while word order between the input and target align in this example, it can differ

https://deeplearning.cs.cmu.edu/F21/document/slides/lec18.attention.pdf

Decoder decides which inputs are needed for prediction at each time step; e.g., "hard attention" focuses on one input



https://deeplearning.cs.cmu.edu/F21/document/slides/lec18.attention.pdf

Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input



Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input



Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input



Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input



"Soft" Attention: Challenge

Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input

Input

He	loved	to	eat

Er liebte zu essen t=1 t=2 t=3 t=4

Target

How should weights be chosen for each input?

"Soft" Attention: Challenge

Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input



Collect manual annotations and then have loss function push predicted parameters to match ground truth parameters... but this is impractical

"Soft" Attention: Challenge

Decoder decides which inputs are needed for prediction at each time step; e.g., "soft attention" uses a weighted combination of the input

Input

He	loved	to	eat

Instead, have the model learn how to weight each input!

Target

Er liebte zu essen t=1 t=2 t=3 t=4



Popular Solutions

A Solution

- 3. At each decoder time step, a prediction is made based on the weighted sum of the inputs
- 2. At each decoder time step,attention weights are computedthat determine each input'srelevance for the prediction

1. Encoder produces hidden state for every input



Today's Topics

- Motivation: machine neural translation for long sentences
- Decoder: attention
- Encoder
- Performance evaluation
- Final project: ways to find a partner

Solution

3. At each decoder time step, a prediction is made based on the weighted sum of the inputs

2. At each decoder time step, attention weights are computed that determine each input's relevance for the prediction



Recall: GRU is a Type of RNN

• Its hidden state captures information about the past



http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/

How many input are considered by the decoder in this example?

At each decoder time step, the compatibility between the decoder's hidden state and each input's hidden state is computed to decide each input's score at the time step

At each decoder time step, the compatibility between the decoder's hidden state and each input's hidden state is computed to decide each input's score at the time step

At each decoder time step, the compatibility between the decoder's hidden state and each input's hidden state is computed to decide each input's score at the time step

At each decoder time step, the compatibility between the decoder's hidden state and each input's hidden state is computed to decide each input's score at the time step

How to measure compatibility between hidden states of the decoder and input?

Many options (function should be differentiable)

Many options (function should be differentiable)

What model parameters must be learned when using dot-product?

Many options (function should be differentiable)

What model parameters must be learned when using bilinear?

Many options (function should be differentiable)

What model parameters must be learned when using multi-layer perceptron?

Many options (function should be differentiable)

Model parameters that must be learned

Next, apply softmax so all inputs' weights sum to 1

We now have our attention weights!

Intuitively:

He loved to eat

Input

The model can weight each input at each time step!

Target Er liebte zu essen

Solution

3. At each decoder time step, a prediction is made based on the weighted sum of the inputs

2. At each decoder time step, attention weights are computed that determine each input's relevance for the prediction

Word Prediction

Compute at time step *t* for all *n* inputs weighted sum:

 $\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \boldsymbol{h}_i$

Influence of inputs are **amplified** for large attention weights and repressed otherwise

Word Prediction

Decoder predicts not only using Decoder its output at the previous time GRU step and previous hidden state, concat but also with a context vector Attention layer addition context vector multiplication multiplication multiplication multiplication alignment vector softmax softmax softmax softmax 1 1 1 decoder score score score score hidden state 0000 encoder $\bigcirc \bigcirc$ $\bigcirc \bigcirc \bigcirc \bigcirc \bigcirc$) hidden state BiGRU Encoder

Bahdanau method

Many options exist for how to use the context vector with the decoder's output at the previous time step to produce an output at each decoder time step

Google method

Many options exist for how to use the context vector with the decoder's output at the previous time step to produce an output at each decoder time step

Decoder $\bigcirc \bigcirc$ $\bigcirc \bigcirc$ $\bigcirc\bigcirc$ LSTM 8 $\bigcirc \bigcirc$ $\bigcirc\bigcirc$ $\bigcirc \bigcirc$ LSTM 4 $\bigcirc\bigcirc$ $\bigcirc \bigcirc$ $\bigcirc\bigcirc$ LSTM 3 $\bigcirc\bigcirc$ LSTM 2 Decoder also uses $\bigcirc \bigcirc$ LSTM 1 previous hidden state concat context vector

Luong method

Many options exist for how to use the context vector with the decoder's output at the previous time step to produce an output at each decoder time step

Decoder

What stays the same at each decoder time step? - input's hidden state

What changes at each decoder time step?

- decoder's hidden state
- (and so) attention weights and context vector
- decoder's prediction at the previous time step

Summary: Attention (Computations at Each Decoder Step)

Decoder decides which inputs are needed for prediction at each time step with "soft attention"

Summary: Attention (Computations at Each Decoder Step)

All parts are differentiable which means end-to-end training is possible

Today's Topics

- Motivation: machine neural translation for long sentences
- Decoder: attention
- Encoder
- Performance evaluation
- Final project: ways to find a partner

Solution

1. Encoder produces hidden state for every input

Recall: Hidden States from RNNs

• Hidden state captures information about the past

http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/

• Two RNNs where input is fed forward and backward respectively and then the hidden states (typically) are concatenated into a hidden state

Input encoded with a bidirectional GRU

• Two RNNs where input is fed forward and backward respectively and then the hidden states (typically) are concatenated into a hidden state

What are advantages of a bi-directional RNN compared to a single RNN?

• Two RNNs where input is fed forward and backward respectively and then the hidden states (typically) are concatenated into a hidden state

Can use information from the past and **future** to make predictions: e.g., can resolve for "Teddy is a ...?" if Teddy refers to a "bear" or former US President Roosevelt

• Two RNNs where input is fed forward and backward respectively and then the hidden states (typically) are concatenated into a hidden state

What are disadvantages of a bi-directional RNN compared to a single RNN?

• Two RNNs where input is fed forward and backward respectively and then the hidden states (typically) are concatenated into a hidden state

Entire sequence must be observed to make a prediction (e.g., unsuitable for text prediction)

Luong's Encoder

Input encoded with a 2-layer stacked LSTM

https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3#df28 Luong et al. Effective Approaches to Attention-based Neural Machine Translation. EMNLP 2015

Google's Encoder

8 layers with a bi-directional first layer and skip connections between layers (greater level of abstraction for input)

https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3#df28

Wu et al. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv 2016.

Today's Topics

- Motivation: machine neural translation for long sentences
- Encoder
- Decoder: attention
- Performance evaluation
- Final project: ways to find a partner

Values are 0 to 1, with whiter pixels indicating larger attention weights

Visualizing Attention: Group Discussion

What insights can we glean from these examples?

While a linear alignment between input and output sentences is common, there are exceptions (e.g., order of adjectives and nouns can differ)

Output words are often informed by more than one input word; e.g., "man" indicates translation of "the" to l' instead of le, la, or les

It naturally handles different input and output lengths (e.g., 1 extra output word for both examples)

Today's Topics

- Motivation: machine neural translation for long sentences
- Encoder
- Decoder: attention
- Performance evaluation
- Final project: ways to find a partner

Final Project

- Requirements on course website: <u>https://dannagurari.colorado.edu/course/neural-networks-and-deep-learning-spring-2025/final-project/</u>
- We will facilitate a 3-stage process to help you find a partner
 - Described in doc linked to from Canvas

Today's Topics

- Motivation: machine neural translation for long sentences
- Encoder
- Decoder: attention
- Performance evaluation
- Final project: ways to find a partner

