# Attention

Tarushii Goel

April 2020

## 1 Overview

Recurrent Neural Networks (RNNs) provide outstanding results on a variety of Natural Language Processing tasks; however, they continue to struggle in capturing long-term dependencies in complex sentences. Several advancements, such as LSTMs & GRUs were designed to tackle the problem of learning long-term dependencies, but perhaps one of the most influential is Attention. This lecture will begin by discussing the origin of Attention as a improvement to Seq2Seq for Neural Machine Translation (NMT), discuss and compare a variety of attention mechanisms, and conclude by exploring self-attention multi-head attention, which serve as the groundwork for recent advancements in Transformers and BERT.
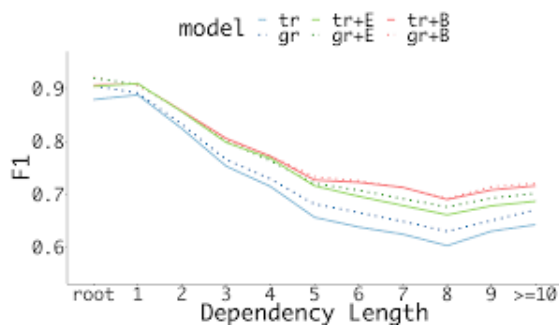


Figure 1: The Long-Term Dependency Problem (Graph of BERT variants, which are state-of-the-art)

## 2 Recap

To refresh you memory, I'll briefly review RNNs and Encoder-Decoders models. RNNs are recurrent; in each iteration they are fed the next input in a sequence and the previous hidden state, which carries accumulated information from other parts of the sentence, and produce a new hidden state.
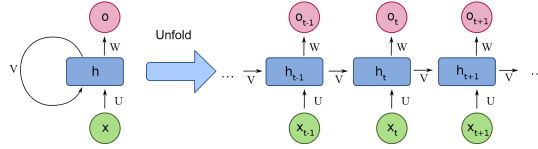
Figure 2: A Diagram of an RNN

Seq2Seq are a unique RNN architecture used for NMT. It is composed to two RNNs; an encoder, which 'encodes' the sentences into a 'context' vector (the last encoder hidden state, a latent representation of the entire sentence), and a decoder, which takes the 'context' vector and decodes it into a tranlated sentence.
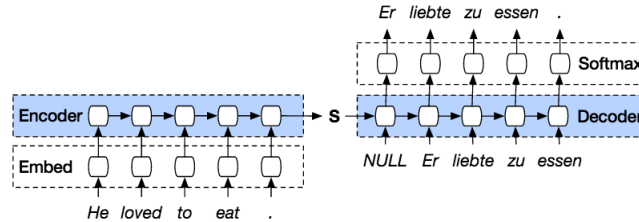


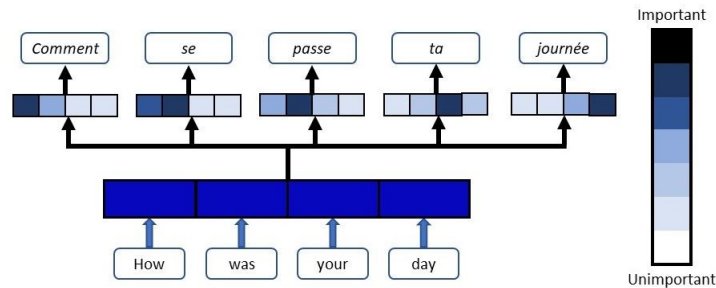Figure 3: A Diagram of an RNN

# 3    What is Attention?

When we think about the English word "Attention", we know that it means directing your focus at something and taking greater notice. The Attention mechanism in Deep Learning is based off this concept of directing your focus, and it pays greater attention to certain factors when processing the data.

In broad terms, Attention is one component of a network's architecture, and is in charge of managing and quantifying the interdependence:

1. Between the input and output elements (General Attention)

2. Within the input elements (Self-Attention)

Say we have the sentence "How was your day", which we would like to translate to the French version - "Comment se passe ta journée". What the Attention component of the network will do for each word in the output sentence is map the important and relevant words from the input sentence and assign higher weights to these words, enhancing the accuracy of the output prediction.

Now, the question remains, how do we calculate attention weights?

Figure 4: Example of Attention

# 4  Types of Attention Calculation

## 4.1  General Attention

General Attention is used to improve Seq2Seq models. Traditionally, the only information carried over from the encoder to the decoder is the last encoder hidden state, but with Attention, a weighted sum of the encoder hidden states of every time-step is appended to the input. Given a query q and a set of key-



Figure 5: General Attention

value pairs (K, V), attention can be generalised to compute a weighted sum of the values dependent on the query and the corresponding keys. The query determines which values to focus on; we can say that the query 'attends' to the values. In the Seq2Seq, the query is the previous decoder hidden state $s_i - 1$ while the set of encoder hidden states $h_0$ to $h_n$ represented both the keys and the

values. The alignment model, is a row of weights applied to the decoder hidden state and the encoder hidden states. To calculate an 'alignment score' for any, calculate the dot-product $s_i^T * h_i$. Once the 'alignment scores are calculated for $h_0 - n$, apply a softmax to get the weights.

## 4.2    Self-Attention

With self-attention, each hidden state attends to the previous hidden states of the same RNN. Here $s_t$ is the query while the decoder hidden states $s_0$ to $s_t - 1$ represent both the keys and the values.
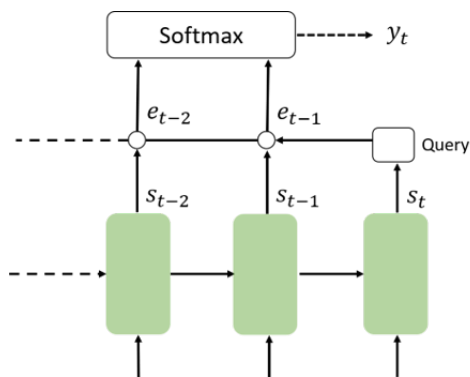


Figure 6: Self-Attention

## 4.3    Multi-Head Attention

When we have multiple queries q, we can stack them in a matrix Q. Each query can be thought of as a different 'head.' In the figure below, two attention heads are displayed (biege and green). The query word is "it". The first attention head is focusing more on the "animal" while the second in green is focusing on "tired". These multi-head mechanisms apparently help in translation tasks.
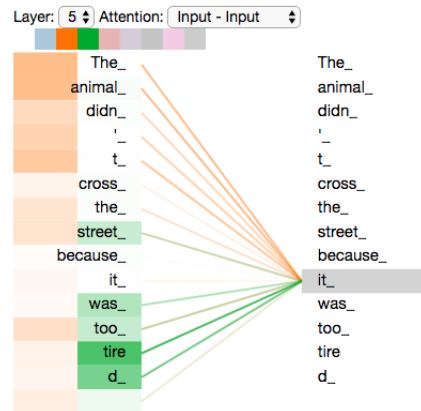
4

Figure 7: Multi-Head Attention

# 5 Code

I coded a standard Seq2Seq model with attention in PyTorch and put the code below. I highly recommend implementing the techniques I discussed in this lecture on your own and using this as a reference. Let me know if you have any questions!

```python
#data is from European Parliament proceedings, download here: http
    ://www.statmt.org/europarl/
PATH = 'YOUR_PATH' #set to the path to your data
en_data = open(PATH + 'europarl-v7.fr-en.en', encoding='utf-8').
    read().split('\n')
fr_data = open(PATH + 'europarl-v7.fr-en.fr', encoding='utf-8').
    read().split('\n')

import spacy
import torch
import torchtext
from torch import nn
import torch.nn.functional as F

device = torch.device("cuda" if torch.cuda.is_available() else "cpu
    ")

en_field = torchtext.data.Field(sequential=True, use_vocab=True,
    init_token='<sos>', eos_token='<eos>', tokenize='spacy',
    tokenizer_language='en')
fr_field = torchtext.data.Field(sequential=True, use_vocab=True,
    init_token = '<sos>', eos_token = '<eos>', tokenize='spacy',
    tokenizer_language='fr')

#Now, covert the data to a csv file to take advantage of torchtext'
    s versatile TabularDataset class
import pandas as pd
from sklearn.model_selection import train_test_split
```

```
21 raw_data = {'English' : [line for line in en_data], 'French': [line
       for line in fr_data]}
22 df = pd.DataFrame(raw_data, columns=["English", "French"])
23 # create train and validation set
24 train, val = train_test_split(df, test_size=0.1)
25 train.to_csv(PATH + "train.csv", index=False)
26 val.to_csv(PATH + "val.csv", index=False)
27
28 train_set,val_set = torchtext.data.TabularDataset.splits(path=PATH,
       train='train.csv', validation='val.csv', format='csv', fields
      =[('English', en_field), ('French', fr_field)])
29
30 #loads 100-dimensional pre-trained glove word embeddings
31 en_field.build_vocab(train_set, val_set, vectors='glove.6B.100d')
32 fr_field.build_vocab(train_set, val_set, vectors='glove.6B.100d')
33
34 train_iter, val_iter = torchtext.data.BucketIterator.splits((
      train_set, val_set), batch_size=32, sort_key=lambda x: len(x.
      French), shuffle=True)
35
36 class Embedder(nn.Module):
37     def __init__(self, embedding):
38         super().__init__()
39         self.embed = nn.Embedding(embedding.shape[0], embedding.
      shape[1])
40         self.embed.weight.data.copy_(embedding)
41         self.embed.weight.requires_grad = False
42     def forward(self, input_sequence):
43         return self.embed(input_sequence)
44
45 class Embedder(nn.Module):
46     def __init__(self, embedding):
47         super().__init__()
48         self.embed = nn.Embedding(embedding.shape[0], embedding.
      shape[1])
49         self.embed.weight.data.copy_(embedding)
50         self.embed.weight.requires_grad = False
51     def forward(self, input_sequence):
52         return self.embed(input_sequence)
53
54 class Encoder(nn.Module):
55     def __init__(self, hidden_size, embedding, num_layers=1,
      dropout=0.0):
56         super(Encoder, self).__init__()
57         self.hidden_size = hidden_size
58         self.embedding = Embedder(embedding)
59         self.gru = nn.GRU(embedding.shape[1], hidden_size,
      num_layers=num_layers, dropout=dropout, bidirectional=True,
      batch_first=True)
60     def forward(self, input_sequence):
61         embedded = self.embedding(input_sequence)
62       #   x = nn.utils.rnn.pack_padded_sequence(x, lens) # unpad
63         outputs, hidden_state = self.gru(embedded) # gru returns
      hidden state of all timesteps as well as hidden state at last
      timestep
64         # pad the sequence to the max length in the batch
65       #   output, _ = nn.utils.rnn.pad_packed_sequence(output)
```

```
66          # The ouput of a GRU has shape (seq_len, batch, hidden_size
        * num_directions)
67          # Because the Encoder is bidirectional, combine the results
         from the
68          # forward and reversed sequence by simply adding them
        together.
69          outputs = outputs[:, :, :self.hidden_size] + outputs[:, :,
        self.hidden_size:]
70          return outputs, hidden_state
71
72 class Decoder(nn.Module):
73      def __init__(self, batch_size, hidden_size, embedding,
        num_layers=1, drop_prob=0.1):
74          super(Decoder, self).__init__()
75          self.batch_size = batch_size
76          self.embedding_size = embedding.shape[1]
77          self.embedding = Embedder(embedding)
78          self.attn = nn.Linear(hidden_size * 2, 1)
79          self.gru = nn.GRU(embedding.shape[1]+hidden_size,
        hidden_size, num_layers=num_layers, dropout=drop_prob,
        batch_first=True)
80          self.classifier = nn.Linear(hidden_size, embedding.shape
        [0])
81
82      def forward(self, decoder_hidden, encoder_outputs, inputs):
83          # Embed input words
84          embedded = self.embedding(inputs)
85          #Assumed size of decoder_hidden -> (num_layers, batch_size,
         embedding_size) size of encoder_outputs -> (batch_size,
        sentence_len, embedding_size)
86          #need to convert length of decoder_hidden to (batch_size,
        sentence_len, embedding_size)
87          self.sequence_len = encoder_outputs.shape[1]
88          decoder_hidden = torch.sum(decoder_hidden, axis=0)
89          attn_inp = decoder_hidden.unsqueeze(1).repeat(1, self.
        sequence_len,1)
90          weights = self.attn(torch.cat((attn_inp, encoder_outputs),
        dim = 2)).squeeze()
91          normalized_weights = F.softmax(weights)
92          attn_applied = torch.bmm(normalized_weights.unsqueeze(1),
        encoder_outputs)
93          cat_input = torch.cat((embedded, attn_applied), axis=2)
94          output, hidden_state = self.gru(cat_input, decoder_hidden.
        unsqueeze(0))
95          output = self.classifier(output).squeeze()
96          return output, hidden_state
97
98 def train(encoder, decoder, encoder_opt, decoder_opt, criterion,
        input, target):
99      #set both to train moode
100     encoder.train()
101     decoder.train()
102     #pass through encoder
103     enc_output, enc_hidden = encoder(target)
104     #initialize input to '<sos>' tokends anddecoder hidden state to
         final encoder hidden state
105     dec_input, dec_hidden = target[:, 0].unsqueeze(1), enc_hidden
```

```
106    loss = 0
107    for i in range(1, target.shape[1]):
108        dec_input, dec_hidden = decoder(dec_hidden, enc_output,
       dec_input)
109        loss += criterion(dec_input, target[:, i])
110        topv, topi = dec_input.topk(1)
111        dec_input = topi.detach()  # detach from history as input
112    loss.backward()
113    encoder_opt.step()
114    decoder_opt.step()
115    return loss
116
117 hidden_size = 200
118 batch_size = 2
119 epochs = 1
120 encoder = Encoder(hidden_size, en_field.vocab.vectors).to(device)
121 decoder = Decoder(batch_size, hidden_size, fr_field.vocab.vectors).
       to(device)
122 encoder_opt = torch.optim.Adam([param for param in encoder.
       parameters() if param.requires_grad == True], lr=1.0e-4)
123 decoder_opt = torch.optim.Adam([param for param in decoder.
       parameters() if param.requires_grad == True], lr=1.0e-4)
124 criterion = nn.CrossEntropyLoss(ignore_index=1)
125
126 for i in range(epochs):
127    losses = []
128    print('Epoch %x:' % i, end='')
129    c = 0
130    for batch in train_iter:
131        input = batch.English.t().to(device)
132        target = batch.French.t().to(device)
133        if (torch.max(input) > 116730):
134            continue
135        losses.append(train(encoder, decoder, encoder_opt,
       decoder_opt, criterion, input, target).item())
136        if (c % 10 == 0):
137            print('.', end='')
138        c+=1
139    print(' loss: ', sum(losses)/len(losses))
```

# 6  References

I got a lot of the explanation and graphics for Attention from this blog post:
https://blog.floydhub.com/attention-mechanism/.