# Sequence-to-Sequence Models and Attention

CS 780/880 Natural Language Processing Lecture 18

Samuel Carton, University of New Hampshire

# Last lecture
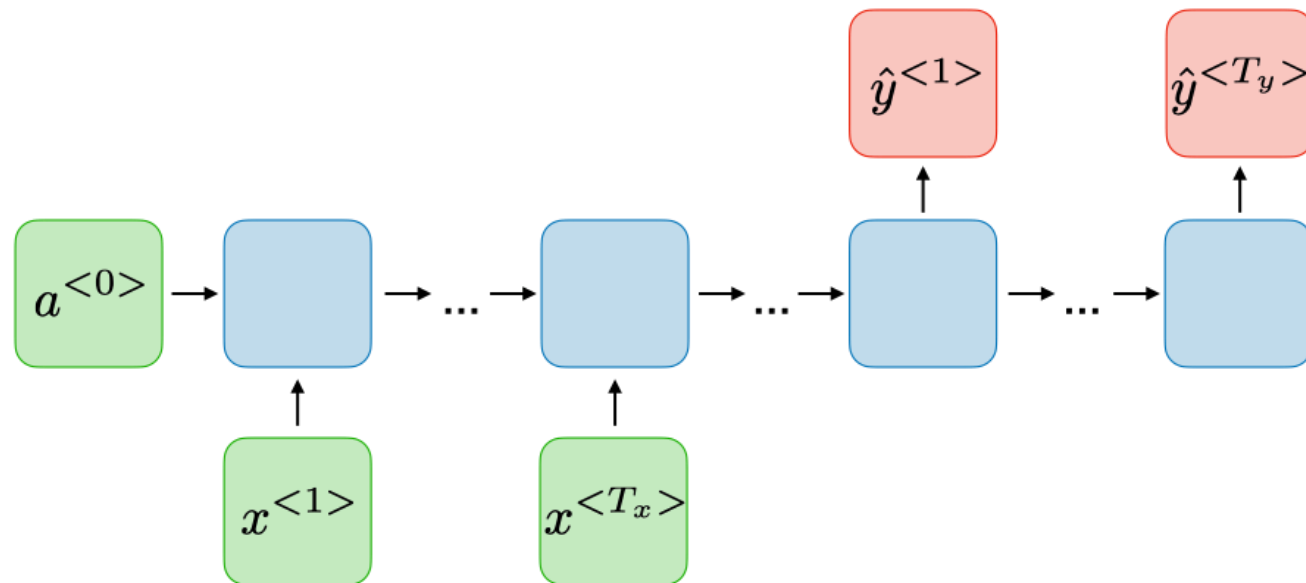
RNNs for language modeling

Generating text
- Greedy decoding
- Random sampling
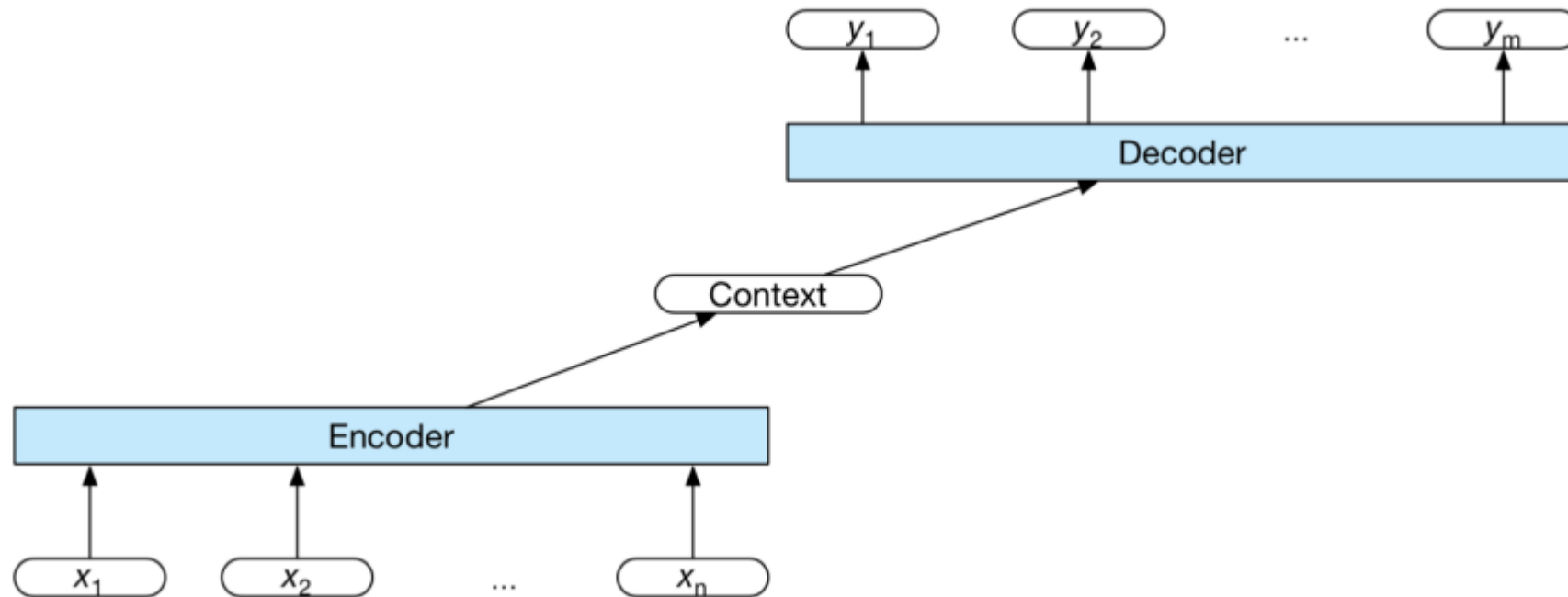- Beam search decoding

Training RNNs
- Teacher forcing
  - Exposure bias
- Alternatives
  - Minimum risk, reinforcement learning, GANs

# Sequence-to-sequence models

**Basic idea**: run an entire sequence through an RNN (the **encoder**), and then give the final vector it makes (the **context**) to another RNN (the **decoder**) to generate a new text sequence with



https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks

3

# Sequence-to-sequence models



https://courses.engr.illinois.edu/cs447/fa2020/Slides/Lecture12.pdf

# Machine translation

**One to-one:**        John loves Mary.
                       |     |     |
                       *Jean aime Marie.*

https://courses.engr.illinois.edu/cs447/fa2020/Slides/Lecture13.pdf

Sequence tagging will work

# Machine translation

**One to-one:**

John loves Mary.

Jean aime Marie.

**One-to-many:**
**(and reordering)**

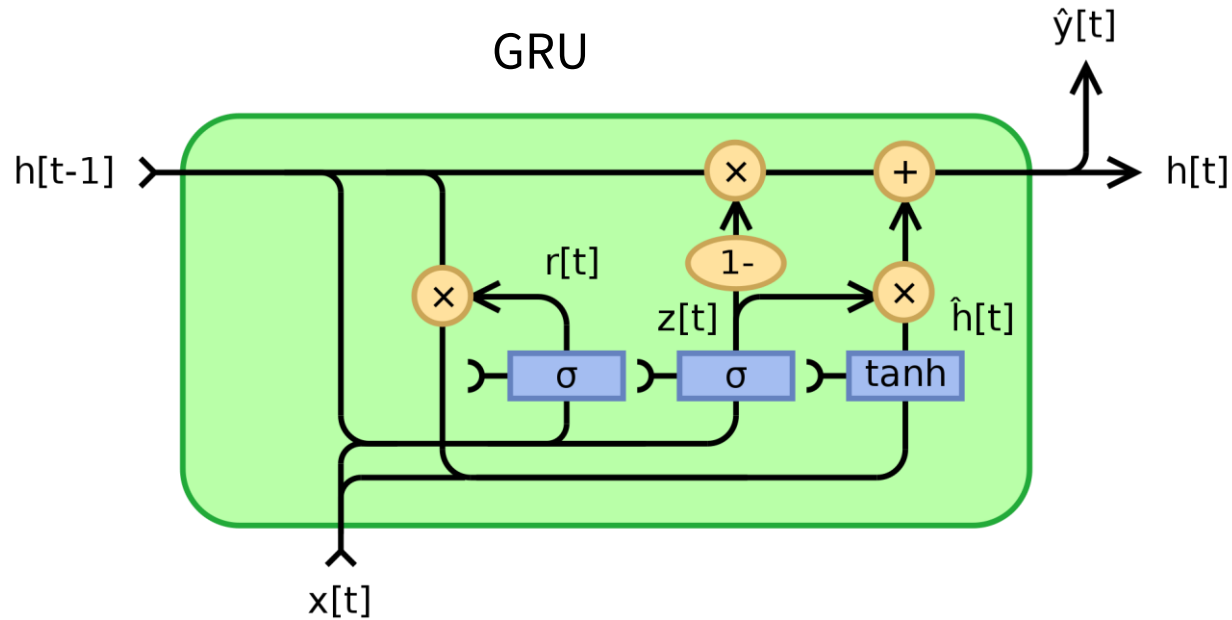John told Mary a story.

Jean [a raconté ] une histoire [à Marie].

**Many-to-one:**
**(and elision)**

John is a [computer scientist].

Jean est informaticien.

**Many-to-many:**

John [swam across] the lake.

Jean [a traversé] le lac [à la nage].

Sequence tagging won't work!

https://courses.engr.illinois.edu/cs447/fa2020/Slides/Lecture13.pdf

# Gated Recurrent Unit (GRU)



GRU

https://en.wikipedia.org/wiki/Gated_recurrent_unit

LSTM

https://en.wikipedia.org/wiki/Long_short-term_memory

# Downloading the translation dataset

```
1 data_url = 'https://download.pytorch.org/tutorial/data.zip'

6 !wget $data_url # this is a linux comand that will grab the file to the local directory
7 !unzip data.zip
```

```
--2023-04-04 16:17:57--  https://download.pytorch.org/tutorial/data.zip
Resolving download.pytorch.org (download.pytorch.org)... 52.222.139.109, 52.222.139.21, 52.222.139.90, ...
Connecting to download.pytorch.org (download.pytorch.org)|52.222.139.109|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 2882130 (2.7M) [application/zip]
Saving to: 'data.zip'

data.zip            100%[===================>]   2.75M  --.-KB/s    in 0.08s

2023-04-04 16:17:58 (34.2 MB/s) - 'data.zip' saved [2882130/2882130]

Archive:  data.zip
   creating: data/
  inflating: data/eng-fra.txt
```

# Downloading the translation dataset

```
1 !ls # this shows us what files we've downloaded
```

data   data.zip   sample_data

```
1 !ls data #this shows us what is inside the file we just unzipped
```

eng-fra.txt   names

# Seq2Seq tutorial

https://pytorch.org/tutorials/intermediate/seq2seq_translation_tutorial.html

A little different from what I've been showing you

- GRU instead of LSTM

- Manual training loop

- Incorporates modeling decisions into training loop

  - I think this is terrible

- Very "from scratch"

  - No word embeddings

  - No NLTK

# Text preprocessing

```python
1  SOS_token = 0
2  EOS_token = 1
3
4
5  class Lang:
6      def __init__(self, name):
7          self.name = name
8          self.word2index = {}
9          self.word2count = {}
10         self.index2word = {0: "SOS", 1: "EOS"}
11         self.n_words = 2  # Count SOS and EOS
12
13     def addSentence(self, sentence):
14         for word in sentence.split(' '):
15             self.addWord(word)
16
17     def addWord(self, word):
18         if word not in self.word2index:
19             self.word2index[word] = self.n_words
20             self.word2count[word] = 1
21             self.index2word[self.n_words] = word
22             self.n_words += 1
23         else:
24             self.word2count[word] += 1
```

```python
1  # Turn a Unicode string to plain ASCII, thanks to
2  # https://stackoverflow.com/a/518232/2809427
3  def unicodeToAscii(s):
4      return ''.join(
5          c for c in unicodedata.normalize('NFD', s)
6          if unicodedata.category(c) != 'Mn'
7      )
8
9  # Lowercase, trim, and remove non-letter characters
10
11
12 def normalizeString(s):
13     s = unicodeToAscii(s.lower().strip())
14     s = re.sub(r"([.!?])", r" \1", s)
15     s = re.sub(r"[^a-zA-Z.!?]+", r" ", s)
16     return s
```

# Text preprocessing

```python
def readLangs(lang1, lang2, reverse=False):
    print("Reading lines...")

    # Read the file and split into lines
    lines = open('data/%s-%s.txt' % (lang1, lang2), encoding='utf-8').\
        read().strip().split('\n')

    # Split every line into pairs and normalize
    pairs = [[normalizeString(s) for s in l.split('\t')] for l in lines]

    # Reverse pairs, make Lang instances
    if reverse:
        pairs = [list(reversed(p)) for p in pairs]
        input_lang = Lang(lang2)
        output_lang = Lang(lang1)
    else:
        input_lang = Lang(lang1)
        output_lang = Lang(lang2)

    return input_lang, output_lang, pairs
```

```python
MAX_LENGTH = 10

eng_prefixes = (
    "i am ", "i m ",
    "he is", "he s ",
    "she is", "she s ",
    "you are", "you re ",
    "we are", "we re ",
    "they are", "they re "
)


def filterPair(p):
    return len(p[0].split(' ')) < MAX_LENGTH and \
        len(p[1].split(' ')) < MAX_LENGTH and \
        p[1].startswith(eng_prefixes)


def filterPairs(pairs):
    return [pair for pair in pairs if filterPair(pair)]
```

# Text preprocessing

```
 1 def prepareData(lang1, lang2, reverse=False):
 2     input_lang, output_lang, pairs = readLangs(lang1, lang2, reverse)
 3     print("Read %s sentence pairs" % len(pairs))
 4     pairs = filterPairs(pairs)
 5     print("Trimmed to %s sentence pairs" % len(pairs))
 6     print("Counting words...")
 7     for pair in pairs:
 8         input_lang.addSentence(pair[0])
 9         output_lang.addSentence(pair[1])
10     print("Counted words:")
11     print(input_lang.name, input_lang.n_words)
12     print(output_lang.name, output_lang.n_words)
13     return input_lang, output_lang, pairs
14
15
16 input_lang, output_lang, pairs = prepareData('eng', 'fra', True)
```

```
Reading lines...
Read 135842 sentence pairs
Trimmed to 10599 sentence pairs
Counting words...
Counted words:
fra 4345
eng 2803
```

```
 1 pprint(pairs[0:5])
```

```
[['j ai ans .', 'i m .'],
 ['je vais bien .', 'i m ok .'],
 ['ca va .', 'i m ok .'],
 ['je suis gras .', 'i m fat .'],
 ['je suis gros .', 'i m fat .']]
```
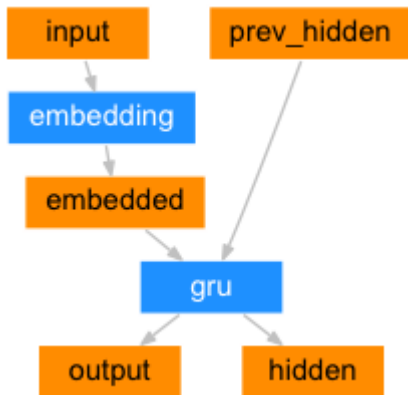
# Model classes

```python
1  class EncoderRNN(nn.Module):
2      def __init__(self, input_size, hidden_size):
3          super(EncoderRNN, self).__init__()
4          self.hidden_size = hidden_size
5
6          self.embedding = nn.Embedding(input_size, hidden_size)
7          self.gru = nn.GRU(hidden_size, hidden_size)
8
9      def forward(self, input, hidden):
10         embedded = self.embedding(input).view(1, 1, -1)
11         output = embedded
12         output, hidden = self.gru(output, hidden)
13         return output, hidden
14     def initHidden(self):
15
16         return torch.zeros(1, 1, self.hidden_size, device=device)
```

```python
1  class DecoderRNN(nn.Module):
2      def __init__(self, hidden_size, output_size):
3          super(DecoderRNN, self).__init__()
4          self.hidden_size = hidden_size
5
6          self.embedding = nn.Embedding(output_size, hidden_size)
7          self.gru = nn.GRU(hidden_size, hidden_size)
8          self.out = nn.Linear(hidden_size, output_size)
9          self.softmax = nn.LogSoftmax(dim=1)
10
11     def forward(self, input, hidden, *args):
12         output = self.embedding(input).view(1, 1, -1)
13         output = F.relu(output)
14         output, hidden = self.gru(output, hidden)
15         output = self.softmax(self.out(output[0]))
16         return output, hidden, None
17
18     def initHidden(self):
19         return torch.zeros(1, 1, self.hidden_size, device=device)
```

# Sequence-to-sequence components

Encoder

Decoder

# Training code

```python
1 def indexesFromSentence(lang, sentence):
2     return [lang.word2index[word] for word in sentence.split(' ')]
3
4
5 def tensorFromSentence(lang, sentence):
6     indexes = indexesFromSentence(lang, sentence)
7     indexes.append(EOS_token)
8     return torch.tensor(indexes, dtype=torch.long, device=device).view(-1, 1)
9
10
11 def tensorsFromPair(pair):
12     input_tensor = tensorFromSentence(input_lang, pair[0])
13     target_tensor = tensorFromSentence(output_lang, pair[1])
14     return (input_tensor, target_tensor)
```

# Training code

```python
1 teacher_forcing_ratio = 0.5
2 def train(input_tensor, target_tensor, encoder, decoder,
3           encoder_optimizer, decoder_optimizer, criterion, max_length=MAX_LENGTH):
4     encoder_hidden = encoder.initHidden()
5
6     encoder_optimizer.zero_grad()
7     decoder_optimizer.zero_grad()
8
9     input_length = input_tensor.size(0)
10    target_length = target_tensor.size(0)
11    encoder_outputs = torch.zeros(max_length, encoder.hidden_size, device=device)
12    loss = 0
13    for ei in range(input_length):
14        encoder_output, encoder_hidden = encoder(
15            input_tensor[ei], encoder_hidden)
16        encoder_outputs[ei] = encoder_output[0, 0]
17
18    decoder_input = torch.tensor([[SOS_token]], device=device)
19    decoder_hidden = encoder_hidden
```

```python
21    use_teacher_forcing = True if random.random() < teacher_forcing_ratio else False
22    if use_teacher_forcing:
23        # Teacher forcing: Feed the target as the next input
24        for di in range(target_length):
25            decoder_output, decoder_hidden, decoder_attention = decoder(
26                decoder_input, decoder_hidden, encoder_outputs)
27            loss += criterion(decoder_output, target_tensor[di])
28            decoder_input = target_tensor[di]  # Teacher forcing
29
30    else:
31        # Without teacher forcing: use its own predictions as the next input
32        for di in range(target_length):
33            decoder_output, decoder_hidden, decoder_attention = decoder(
34                decoder_input, decoder_hidden, encoder_outputs)
35            topv, topi = decoder_output.topk(1)
36            decoder_input = topi.squeeze().detach()  # detach from history as input
37
38            loss += criterion(decoder_output, target_tensor[di])
39            if decoder_input.item() == EOS_token:
40                break
41
42    loss.backward()
43    encoder_optimizer.step()
44    decoder_optimizer.step()
45    return loss.item() / target_length
```

# Training code

```python
1 teacher_forcing_ratio = 0.5
2 def train(input_tensor, target_tensor, encoder, decoder,
3           encoder_optimizer, decoder_optimizer, criterion, max_length=MAX_LENGTH):
4     encoder_hidden = encoder.initHidden()
5
6     encoder_optimizer.zero_grad()
7     decoder_optimizer.zero_grad()
8
9     input_length = input_tensor.size(0)
10    target_length = target_tensor.size(0)
11    encoder_outputs = torch.zeros(max_length, encoder.hidden_size, device=device)
12    loss = 0
13    for ei in range(input_length):
14        encoder_output, encoder_hidden = encoder(
15            input_tensor[ei], encoder_hidden)
16        encoder_outputs[ei] = encoder_output[0, 0]
17
18    decoder_input = torch.tensor([[SOS_token]], device=device)
19    decoder_hidden = encoder_hidden
```

```python
21    use_teacher_forcing = True if random.random() < teacher_forcing_ratio else False
22    if use_teacher_forcing:
23        # Teacher forcing: Feed the target as the next input
24        for di in range(target_length):
25            decoder_output, decoder_hidden, decoder_attention = decoder(
26                decoder_input, decoder_hidden, encoder_outputs)
27            loss += criterion(decoder_output, target_tensor[di])
28            decoder_input = target_tensor[di]  # Teacher forcing
29
30    else:
31        # Without teacher forcing: use its own predictions as the next input
32        for di in range(target_length):
33            decoder_output, decoder_hidden, decoder_attention = decoder(
34                decoder_input, decoder_hidden, encoder_outputs)
35            topv, topi = decoder_output.topk(1)
36            decoder_input = topi.squeeze().detach()  # detach from history as input
37
38            loss += criterion(decoder_output, target_tensor[di])
39            if decoder_input.item() == EOS_token:
40                break
41
42    loss.backward()
43    encoder_optimizer.step()
44    decoder_optimizer.step()
45    return loss.item() / target_length
```

# Training code

```python
1  teacher_forcing_ratio = 0.5
2  def train(input_tensor, target_tensor, encoder, decoder,
3           encoder_optimizer, decoder_optimizer, criterion, max_length=MAX_LENGTH):
4      encoder_hidden = encoder.initHidden()
5
6      encoder_optimizer.zero_grad()
7      decoder_optimizer.zero_grad()
8
9      input_length = input_tensor.size(0)
10     target_length = target_tensor.size(0)
11     encoder_outputs = torch.zeros(max_length, encoder.hidden_size, device=device)
12     loss = 0
13     for ei in range(input_length):
14         encoder_output, encoder_hidden = encoder(
15             input_tensor[ei], encoder_hidden)
16         encoder_outputs[ei] = encoder_output[0, 0]
17
18     decoder_input = torch.tensor([[SOS_token]], device=device)
19     decoder_hidden = encoder_hidden
```

```python
21     use_teacher_forcing = True if random.random() < teacher_forcing_ratio else False
22     if use_teacher_forcing:
23         # Teacher forcing: Feed the target as the next input
24         for di in range(target_length):
25             decoder_output, decoder_hidden, decoder_attention = decoder(
26                 decoder_input, decoder_hidden, encoder_outputs)
27             loss += criterion(decoder_output, target_tensor[di])
28             decoder_input = target_tensor[di]  # Teacher forcing
29
30     else:
31         # Without teacher forcing: use its own predictions as the next input
32         for di in range(target_length):
33             decoder_output, decoder_hidden, decoder_attention = decoder(
34                 decoder_input, decoder_hidden, encoder_outputs)
35             topv, topi = decoder_output.topk(1)
36             decoder_input = topi.squeeze().detach()  # detach from history as input
37
38             loss += criterion(decoder_output, target_tensor[di])
39             if decoder_input.item() == EOS_token:
40                 break
41
42     loss.backward()
43     encoder_optimizer.step()
44     decoder_optimizer.step()
45     return loss.item() / target_length
```

# Training code

```python
1 teacher_forcing_ratio = 0.5
2 def train(input_tensor, target_tensor, encoder, decoder,
3           encoder_optimizer, decoder_optimizer, criterion, max_length=MAX_LENGTH):
4     encoder_hidden = encoder.initHidden()
5
6     encoder_optimizer.zero_grad()
7     decoder_optimizer.zero_grad()
8
9     input_length = input_tensor.size(0)
10    target_length = target_tensor.size(0)
11    encoder_outputs = torch.zeros(max_length, encoder.hidden_size, device=device)
12    loss = 0
13    for ei in range(input_length):
14        encoder_output, encoder_hidden = encoder(
15            input_tensor[ei], encoder_hidden)
16        encoder_outputs[ei] = encoder_output[0, 0]
17
18    decoder_input = torch.tensor([[SOS_token]], device=device)
19    decoder_hidden = encoder_hidden
```

```python
21    use_teacher_forcing = True if random.random() < teacher_forcing_ratio else False
22    if use_teacher_forcing:
23        # Teacher forcing: Feed the target as the next input
24        for di in range(target_length):
25            decoder_output, decoder_hidden, decoder_attention = decoder(
26                decoder_input, decoder_hidden, encoder_outputs)
27            loss += criterion(decoder_output, target_tensor[di])
28            decoder_input = target_tensor[di]  # Teacher forcing
29
30    else:
31        # Without teacher forcing: use its own predictions as the next input
32        for di in range(target_length):
33            decoder_output, decoder_hidden, decoder_attention = decoder(
34                decoder_input, decoder_hidden, encoder_outputs)
35            topv, topi = decoder_output.topk(1)
36            decoder_input = topi.squeeze().detach()  # detach from history as input
37
38            loss += criterion(decoder_output, target_tensor[di])
39            if decoder_input.item() == EOS_token:
40                break
41
42    loss.backward()
43    encoder_optimizer.step()
44    decoder_optimizer.step()
45    return loss.item() / target_length
```

# Training code

```python
 1 def trainIters(encoder, decoder, n_iters, print_every=1000, plot_every=100, learning_rate=0.01):
 2     start = time.time()
 3     plot_losses = []
 4     print_loss_total = 0  # Reset every print_every
 5     plot_loss_total = 0  # Reset every plot_every
 6
 7     encoder_optimizer = optim.SGD(encoder.parameters(), lr=learning_rate)
 8     decoder_optimizer = optim.SGD(decoder.parameters(), lr=learning_rate)
 9     training_pairs = [tensorsFromPair(random.choice(pairs))
10                       for i in range(n_iters)]
11     criterion = nn.NLLLoss()
12
13     for iter in range(1, n_iters + 1):
14         training_pair = training_pairs[iter - 1]
15         input_tensor = training_pair[0]
16         target_tensor = training_pair[1]
17
18         loss = train(input_tensor, target_tensor, encoder,
19                      decoder, encoder_optimizer, decoder_optimizer, criterion)
20         print_loss_total += loss
21         plot_loss_total += loss
22
23         if iter % print_every == 0:
24             print_loss_avg = print_loss_total / print_every
25             print_loss_total = 0
26             print('%s (%d %d%%) %.4f' % (timeSince(start, iter / n_iters),
27                                          iter, iter / n_iters * 100, print_loss_avg))
28
29         if iter % plot_every == 0:
30             plot_loss_avg = plot_loss_total / plot_every
31             plot_losses.append(plot_loss_avg)
32             plot_loss_total = 0
33
34     showPlot(plot_losses)
```

# Training code

```python
 1 def trainIters(encoder, decoder, n_iters, print_every=1000, plot_every=100, learning_rate=0.01):
 2     start = time.time()
 3     plot_losses = []
 4     print_loss_total = 0  # Reset every print_every
 5     plot_loss_total = 0  # Reset every plot_every
 6
 7     encoder_optimizer = optim.SGD(encoder.parameters(), lr=learning_rate)
 8     decoder_optimizer = optim.SGD(decoder.parameters(), lr=learning_rate)
 9     training_pairs = [tensorsFromPair(random.choice(pairs))
10                       for i in range(n_iters)]
11     criterion = nn.NLLLoss()
12
13     for iter in range(1, n_iters + 1):
14         training_pair = training_pairs[iter - 1]
15         input_tensor = training_pair[0]
16         target_tensor = training_pair[1]
17
18         loss = train(input_tensor, target_tensor, encoder,
19                      decoder, encoder_optimizer, decoder_optimizer, criterion)
20         print_loss_total += loss
21         plot_loss_total += loss
22
23         if iter % print_every == 0:
24             print_loss_avg = print_loss_total / print_every
25             print_loss_total = 0
26             print('%s (%d %d%%) %.4f' % (timeSince(start, iter / n_iters),
27                                          iter, iter / n_iters * 100, print_loss_avg))
28
29         if iter % plot_every == 0:
30             plot_loss_avg = plot_loss_total / plot_every
31             plot_losses.append(plot_loss_avg)
32             plot_loss_total = 0
33
34     showPlot(plot_losses)
```

22

# Model training

```
1 # First we'll try training the version that doesn't use attention
2 hidden_size = 256
3 encoder1 = EncoderRNN(input_lang.n_words, hidden_size).to(device)
4 nonattention_decoder = DecoderRNN(hidden_size, output_lang.n_words).to(device)
5
6 trainIters(encoder1, nonattention_decoder, 25000, print_every=5000)
```

```
0m 56s (- 3m 44s) (5000 20%) 2.9316
1m 49s (- 2m 43s) (10000 40%) 2.3708
2m 42s (- 1m 48s) (15000 60%) 2.0766
3m 36s (- 0m 54s) (20000 80%) 1.8224
4m 30s (- 0m 0s) (25000 100%) 1.6010
```

# Improving naïve seq2seq

**Big problem here**: we're expecting a **lot** out of that final encoder context vector.

- Essentially we're asking it to save up everything it needs to know to then go ahead and spit out the text we want.
- That's a lot of info to squeeze into a 100-element vector

**Idea**: What if we also let the decoder look at the original input while it is decoding the context?

- But it would need to be able to learn which parts of the original input were pertinent to what it was trying to do at any given point

Solution: **Attention**

# Attention in sequence-to-sequence

**Basic idea**: The decoder will have access to all the output from the encoder (not just the final output), but will learn a **weighting function** for how important any individual output is at a given timestep.

Encoder

Decoder (no attention)

Decoder (with attention)

# Model classes

```python
class DecoderRNN(nn.Module):
    def __init__(self, hidden_size, output_size):
        super(DecoderRNN, self).__init__()
        self.hidden_size = hidden_size

        self.embedding = nn.Embedding(output_size, hidden_size)
        self.gru = nn.GRU(hidden_size, hidden_size)
        self.out = nn.Linear(hidden_size, output_size)
        self.softmax = nn.LogSoftmax(dim=1)

    def forward(self, input, hidden, *args):
        output = self.embedding(input).view(1, 1, -1)
        output = F.relu(output)
        output, hidden = self.gru(output, hidden)
        output = self.softmax(self.out(output[0]))
        return output, hidden, None

    def initHidden(self):
        return torch.zeros(1, 1, self.hidden_size, device=device)
```
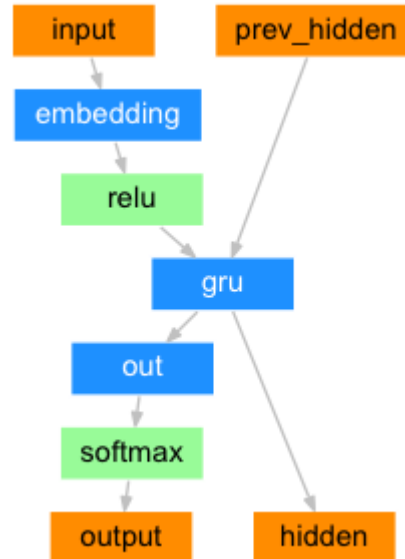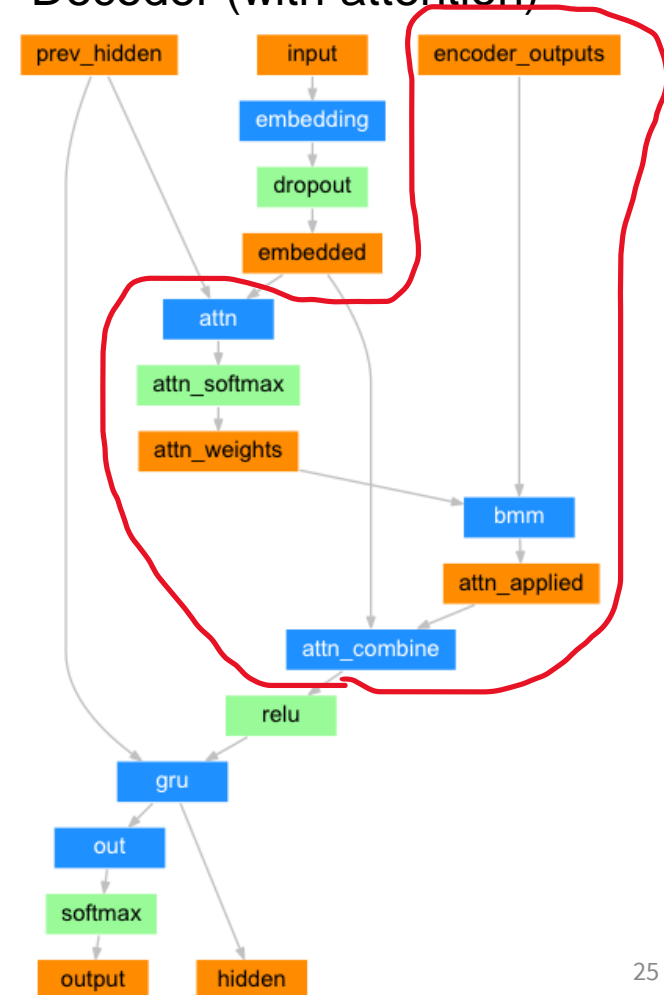
```python
class AttnDecoderRNN(nn.Module):
    def __init__(self, hidden_size, output_size, dropout_p=0.1, max_length=MAX_LENGTH):
        super(AttnDecoderRNN, self).__init__()
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.dropout_p = dropout_p
        self.max_length = max_length

        self.embedding = nn.Embedding(self.output_size, self.hidden_size)
        self.attn = nn.Linear(self.hidden_size * 2, self.max_length)
        self.attn_combine = nn.Linear(self.hidden_size * 2, self.hidden_size)
        self.dropout = nn.Dropout(self.dropout_p)
        self.gru = nn.GRU(self.hidden_size, self.hidden_size)
        self.out = nn.Linear(self.hidden_size, self.output_size)

    def forward(self, input, hidden, encoder_outputs):
        embedded = self.embedding(input).view(1, 1, -1)
        embedded = self.dropout(embedded)

        attn_weights = F.softmax(
            self.attn(torch.cat((embedded[0], hidden[0]), 1)), dim=1)
        attn_applied = torch.bmm(attn_weights.unsqueeze(0),
                                 encoder_outputs.unsqueeze(0))
        output = torch.cat((embedded[0], attn_applied[0]), 1)
        output = self.attn_combine(output).unsqueeze(0)
        output = F.relu(output)
        output, hidden = self.gru(output, hidden)
        output = F.log_softmax(self.out(output[0]), dim=1)
        return output, hidden, attn_weights

    def initHidden(self):
        return torch.zeros(1, 1, self.hidden_size, device=device)
```

# Model training

**Without attention**

```
2 hidden_size = 256
3 encoder1 = EncoderRNN(input_lang.n_words, hidden_size).to(device)
4 nonattention_decoder = DecoderRNN(hidden_size, output_lang.n_words).to(device)
5
6 trainIters(encoder1, nonattention_decoder, 25000, print_every=5000)
```

```
0m 56s (- 3m 44s) (5000 20%) 2.9316
1m 49s (- 2m 43s) (10000 40%) 2.3708
2m 42s (- 1m 48s) (15000 60%) 2.0766
3m 36s (- 0m 54s) (20000 80%) 1.8224
4m 30s (- 0m 0s) (25000 100%) 1.6010
```

**With attention**

```
2 hidden_size = 256
3 encoder2 = EncoderRNN(input_lang.n_words, hidden_size).to(device)
4 attention_decoder = AttnDecoderRNN(hidden_size, output_lang.n_words,
5                                    dropout_p=0.1).to(device)
6
7 trainIters(encoder2, attention_decoder, 25000, print_every=5000)
```

```
1m 23s (- 5m 32s) (5000 20%) 2.8722
2m 48s (- 4m 13s) (10000 40%) 2.2846
4m 15s (- 2m 50s) (15000 60%) 1.9739
5m 38s (- 1m 24s) (20000 80%) 1.7127
7m 1s (- 0m 0s) (25000 100%) 1.5260
```

# Classification with attention

**Basic idea**: Use one RNN (attender) to generate attention weights over a sequence, then a second RNN (predictor) to make predictions from the attention-weighted sequence

Dual training objective which encourages attention weights to be sparse, but predictor to be accurate.

In theory, leads to only important information (stuff needed for prediction) to be attended to.

Prediction

| Predictor |

You are a **huge jerk**

**0    0    0    1    1**

| Attender |

You are a huge jerk

# Attention classification model

```python
class AttentionClassifier(pl.LightningModule):
  def __init__(self,
               word_vectors:np.ndarray,
               num_classes:int,
               learning_rate:float,
               padding_id:int,
               lstm_hidden_size:int=100, # how big the inner vectors of the LSTM will be,
               lstm_layers:int =2, # how many layers the LSTM will have
               dropout_prob:float=0.1,
               sparsity_loss_weight:float= 0.15,
               **kwargs):
    super().__init__( **kwargs)

    # We'll use the same PyTorch Embedding layer as before
    self.word_embeddings = torch.nn.Embedding.from_pretrained(embeddings=torch.tensor(word_vectors),
                                                              freeze=True)


    self.attender = torch.nn.LSTM(input_size = word_vectors.shape[1],
                                  hidden_size = lstm_hidden_size,
                                  num_layers=lstm_layers,
                                  bidirectional=True,
                                  dropout=dropout_prob,
                                  batch_first=True)
    self.attender_output_layer  = torch.nn.Linear(2*lstm_hidden_size, 1)
```

```python
    self.predictor = torch.nn.LSTM(input_size = word_vectors.shape[1],
                                   hidden_size = lstm_hidden_size,
                                   num_layers=lstm_layers,
                                   bidirectional=True,
                                   dropout=dropout_prob,
                                   batch_first=True)
    self.predictor_output_layer  = torch.nn.Linear(2*lstm_hidden_size, num_classes)


    # Output layer input size has to be doubled because the LSTM is bidirectional
    self.lstm_layers = lstm_layers
    self.learning_rate = learning_rate
    self.padding_id = padding_id # we'll need this later
    self.sparsity_loss_weight = sparsity_loss_weight
    self.train_accuracy = Accuracy(task='multiclass', num_classes=num_classes)
    self.val_accuracy = Accuracy(task='multiclass', num_classes=num_classes)
```

# Attention classification model

```python
44  def forward(self, y:torch.Tensor, input_ids:torch.Tensor, verbose=False):
45      inputs_embeds = self.word_embeddings(input_ids) #(batch size x sequence length x embedding size)
46      input_lengths = (input_ids != self.padding_id).sum(dim=1).detach().cpu()
47
48      packed_embeddings = pack_padded_sequence(inputs_embeds, input_lengths, batch_first=True, enforce_sorted=False)
49      packed_attender_output, _ = self.attender.forward(packed_embeddings)
50      attender_output, _ = pad_packed_sequence(packed_attender_output, batch_first=True, padding_value=0.0, total_length=input_ids.shape[1])
51      attention_logits = self.attender_output_layer(attender_output) #(batch size x sequence length x 1)
52      attention_mask = torch.nn.functional.sigmoid(attention_logits)
53      attention_masked_inputs_embeds = attention_mask * inputs_embeds
54      attention_mask = attention_mask.squeeze(-1)
55      sparsity_loss = masked_mean(attention_mask, (input_ids == self.padding_id)).mean()
56
57      packed_masked_embeddings = pack_padded_sequence(attention_masked_inputs_embeds, input_lengths, batch_first=True, enforce_sorted=False)
58      _, (final_predictor_hidden, final_predictor_state) = self.attender.forward(packed_masked_embeddings)
59      last_layer_idx = self.lstm_layers-1
60      last_layer_final_forward_hiddens = final_predictor_hidden[2*last_layer_idx]
61      last_layer_final_reverse_hiddens = final_predictor_hidden[2*last_layer_idx+1]
62      combined_last_layer_hiddens = torch.cat([last_layer_final_forward_hiddens, last_layer_final_reverse_hiddens], dim=1)
63      py_logits = self.predictor_output_layer(combined_last_layer_hiddens)
64      py = torch.argmax(py_logits, dim=1)
65      py_loss = torch.nn.functional.cross_entropy(py_logits, y, reduction='mean')
66
67      loss = py_loss + self.sparsity_loss_weight * sparsity_loss
68      return {'py':py,
69              'sparsity_loss':sparsity_loss,
70              'py_loss':py_loss,
71              'attention_mask':attention_mask,
72              'loss':loss}
```

# Trainer

```python
1 from pytorch_lightning import Trainer
2 from pytorch_lightning.callbacks.progress import TQDMProgressBar
3 from pytorch_lightning.callbacks import ModelCheckpoint
4
5 checkpoint_callback = ModelCheckpoint(dirpath=".", save_top_k=1, monitor="val_loss")
6
7 trainer = Trainer(
8     accelerator="auto",
9     devices=1 if torch.cuda.is_available() else None,
10    max_epochs=3,
11    callbacks=[TQDMProgressBar(refresh_rate=20), checkpoint_callback],
12    val_check_interval = 0.5,
13    default_root_dir='.' # This tells Pytorch Lightning to save checkpoints in the current working directory
14    )
```

```
INFO:pytorch_lightning.utilities.rank_zero:GPU available: True (cuda), used: True
INFO:pytorch_lightning.utilities.rank_zero:TPU available: False, using: 0 TPU cores
INFO:pytorch_lightning.utilities.rank_zero:IPU available: False, using: 0 IPUs
INFO:pytorch_lightning.utilities.rank_zero:HPU available: False, using: 0 HPUs
```

# Trainer

```
1 trainer.fit(model=model,
2              train_dataloaders=train_dataloader,
3              val_dataloaders=dev_dataloader)
```

```
/usr/local/lib/python3.9/dist-packages/pytorch_lightning/callbacks/model_checkpoint.py:613: UserWarning: Checkpoint directory /content exists and is not empty.
  rank_zero_warn(f"Checkpoint directory {dirpath} exists and is not empty.")
INFO:pytorch_lightning.accelerators.cuda:LOCAL_RANK: 0 - CUDA_VISIBLE_DEVICES: [0]
INFO:pytorch_lightning.callbacks.model_summary:
  | Name                  | Type             | Params
-------------------------------------------------------------
0 | word_embeddings       | Embedding        | 40.0 M
1 | attender              | LSTM             | 403 K
2 | attender_output_layer | Linear           | 201
3 | predictor             | LSTM             | 403 K
4 | predictor_output_layer| Linear           | 402
5 | train_accuracy        | MulticlassAccuracy | 0
6 | val_accuracy          | MulticlassAccuracy | 0
-------------------------------------------------------------
807 K      Trainable params
40.0 M     Non-trainable params
40.8 M     Total params
163.229    Total estimated model params size (MB)
Validation accuracy: tensor(0.5234, device='cuda:0')
```

Epoch 2: 100% ████████████████████████████ 1081/1081 [00:24<00:00, 43.77it/s, loss=0.256, v_num=6]

```
Validation accuracy: tensor(0.8016, device='cuda:0')
Validation accuracy: tensor(0.8062, device='cuda:0')
Training accuracy: tensor(0.8196, device='cuda:0')
Validation accuracy: tensor(0.8417, device='cuda:0')
Validation accuracy: tensor(0.8394, device='cuda:0')
Training accuracy: tensor(0.8719, device='cuda:0')
Validation accuracy: tensor(0.8326, device='cuda:0')
INFO:pytorch_lightning.utilities.rank_zero:`Trainer.fit` stopped: `max_epochs=3` reached.
Validation accuracy: tensor(0.8532, device='cuda:0')
Training accuracy: tensor(0.8981, device='cuda:0')
```

# Visualizing model output

```python
1 sentence = "It was a horrible movie, quite literally the most disgusting thing I have ever seen."
2 sentence_label = 0
3
4 tokens = tokenize(sentence)
5 word_ids = tokens_to_ids(tokens)
6
7 input_ids = torch.tensor([word_ids])
8 print(input_ids)
9
10 y = torch.tensor([sentence_label])
11 print(y)
```

```
tensor([[    20,     15,      7, 10230,   1005,      1,   1689,   5917,      0,     96,
          23967,    873,     41,     33,    661,    541,      2]])
tensor([0])
```

```python
1 with torch.no_grad():
2   model_output = model.forward(input_ids=input_ids, y=y)
3 pprint(model_output)
```

```
{'attention_mask': tensor([[0.2119, 0.9403, 0.3098, 0.9849, 0.2164, 0.2423, 0.8949, 0.9950, 0.1231,
         0.0810, 0.6699, 0.3581, 0.4085, 0.2120, 0.9551, 0.1834, 0.0361]]),
 'loss': tensor(0.0012),
 'py': tensor([0]),
 'py_loss': tensor(0.0012),
 'sparsity_loss': tensor(0.)}
```

# Visualizing model output

```
1  from IPython.core.display import HTML
2
3  for token, attention_weight in zip(tokens, model_output['attention_mask'][0]):
4      # print(token, attention_weight)
5      token_html = HTML(f'<span style="background-color: rgba(255,0,0, {attention_weight});">{token}</span>')
6      display(token_html,)
```

it
was
a
horrible
movie
,
quite
literally
the
most
disgusting
thing
i
have
ever
seen
.

# Saving and loading the model

```
1 # We can see the best checkpoint that Pytorch lightning saved for us
2 !ls
```

```
data    data.zip   'epoch=1-step=2105.ckpt'    lightning_logs    sample_data
```

```
1 # But we can also manually save the model in its current state
2 torch.save(model.state_dict(), 'manually_saved_model.ckpt')
```

```
1 !ls
```

```
data       'epoch=1-step=2105.ckpt'    manually_saved_model.ckpt
data.zip    lightning_logs             sample_data
```

# Saving and loading the model

```
 1 loaded_model = AttentionClassifier(word_vectors=vector_model.vectors,
 2                                    num_classes = 2,
 3                                    learning_rate = 0.001, #I'll typically start with something like 1e-3 for LSTMs
 4                                    padding_id = vector_model.key_to_index['<pad>'],
 5                                    lstm_hidden_size=100,
 6                                    lstm_layers=2,
 7                                    dropout_prob=0.1,
 8                                    sparsity_loss_weight=0.15)
 9 loaded_model.load_state_dict(torch.load('manually_saved_model.ckpt'))
10 display(loaded_model)
```

```
AttentionClassifier(
  (word_embeddings): Embedding(400002, 100)
  (attender): LSTM(100, 100, num_layers=2, batch_first=True, dropout=0.1, bidirectional=True)
  (attender_output_layer): Linear(in_features=200, out_features=1, bias=True)
  (predictor): LSTM(100, 100, num_layers=2, batch_first=True, dropout=0.1, bidirectional=True)
  (predictor_output_layer): Linear(in_features=200, out_features=2, bias=True)
  (train_accuracy): MulticlassAccuracy()
  (val_accuracy): MulticlassAccuracy()
)
```

# Concluding thoughts

Sequence-to-sequence models

- Main application: translation

Attention

- Improves performance of sequence-to-sequence models
- Improves interpretability of classifiers

Model saving/loading