



Sequence-to-Sequence Models and Basic Attention

CS 780/880 Natural Language Processing Lecture 19

Samuel Carton, University of New Hampshire



Last lecture

RNNs for language modeling in PyTorch

Generating text

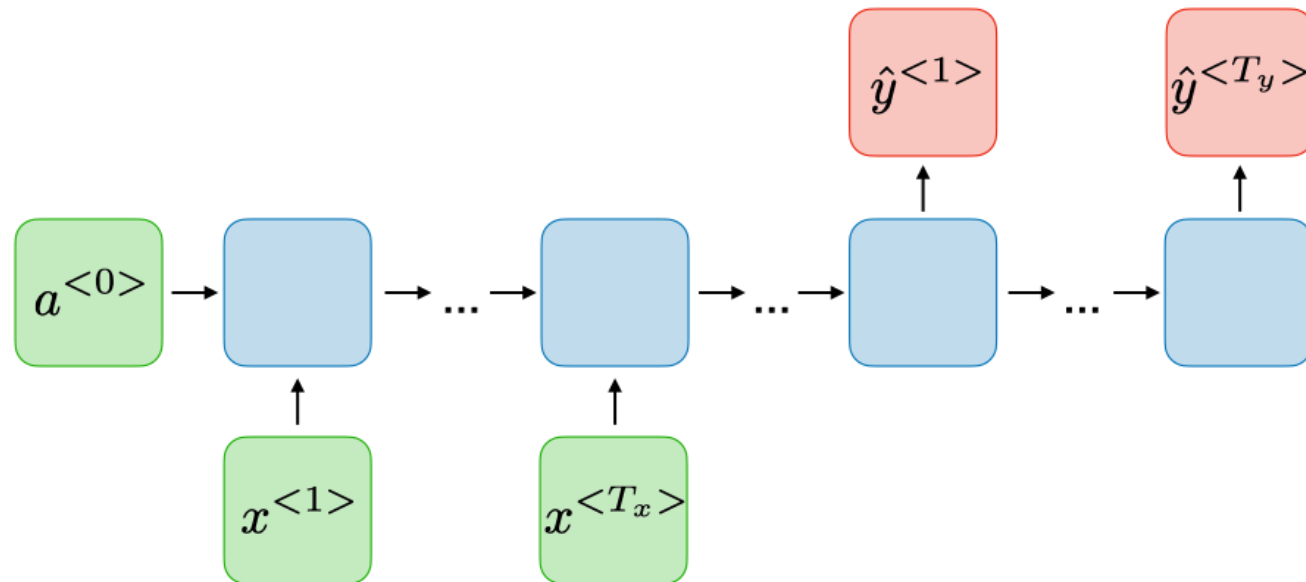
- **Greedy decoding**
- **Random sampling**
- Beam search decoding

Training RNNs

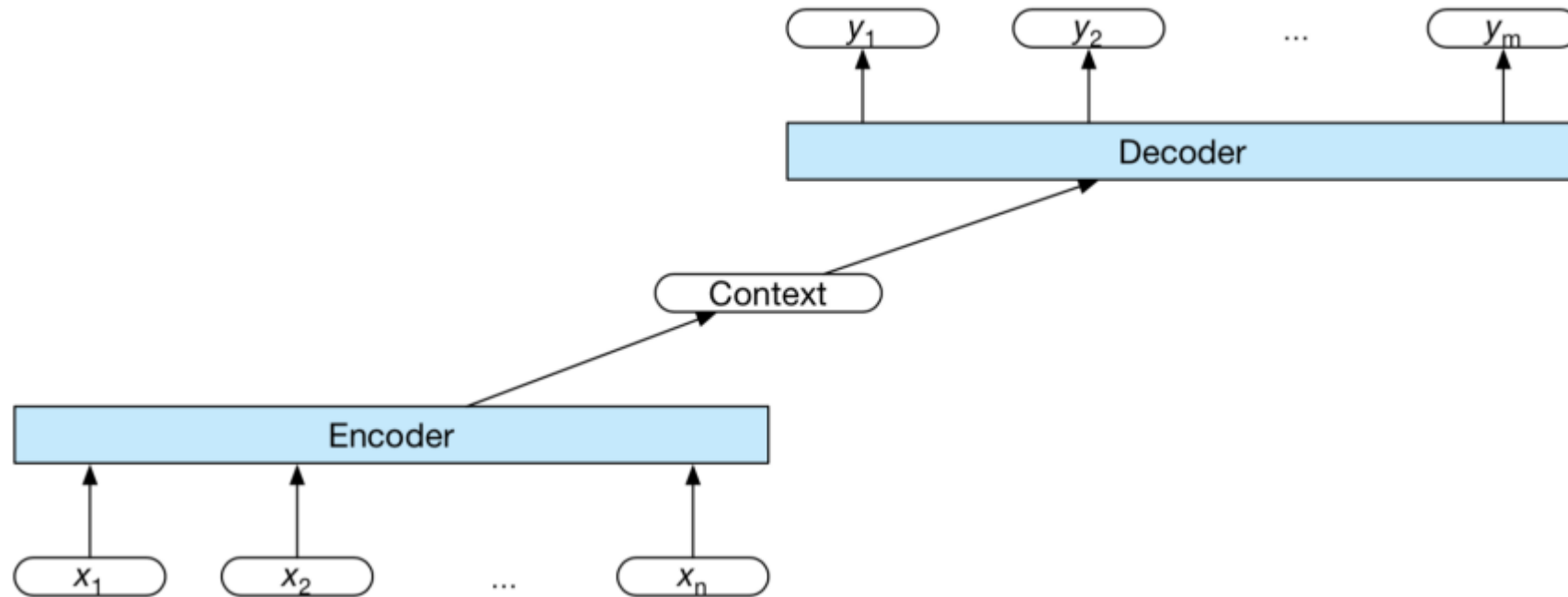
- **Teacher forcing**
 - Exposure bias
- Alternatives
 - Minimum risk, reinforcement learning, GANs

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



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.

Sequence tagging will work

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



Machine translation

One to-one:

John loves Mary.
| | |
Jean aime Marie.

Sequence tagging won't work!

**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:

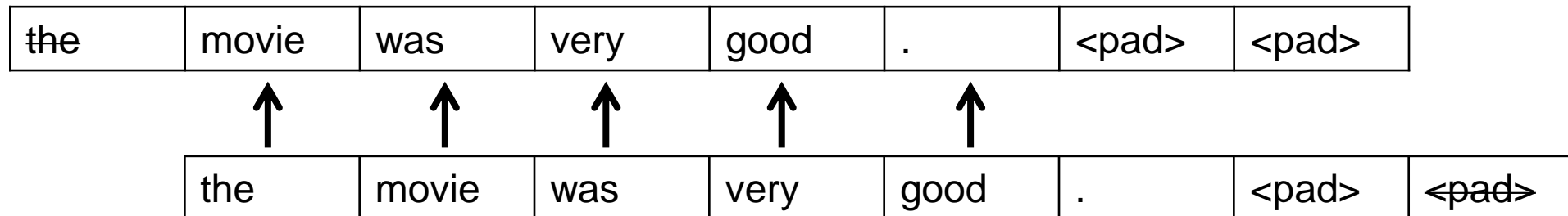
John [swam across] the lake.
| | | |
Jean [a traversé] le lac [à la nage].



Language modeling batch loss

When we do teacher forcing for language modeling, we judge the model's output for each input token, against the next consecutive output token.

But that won't work directly for translation. So what do we do?

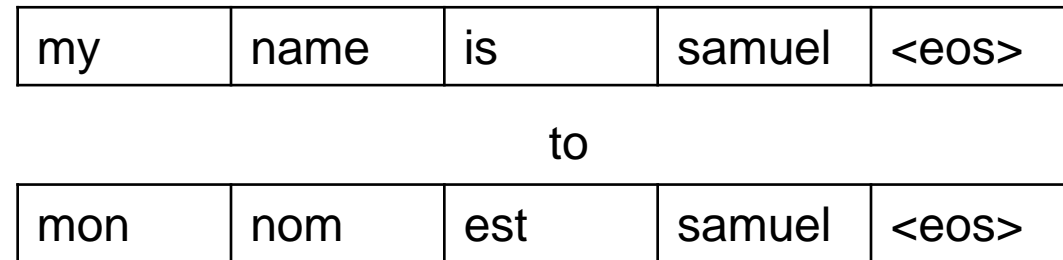




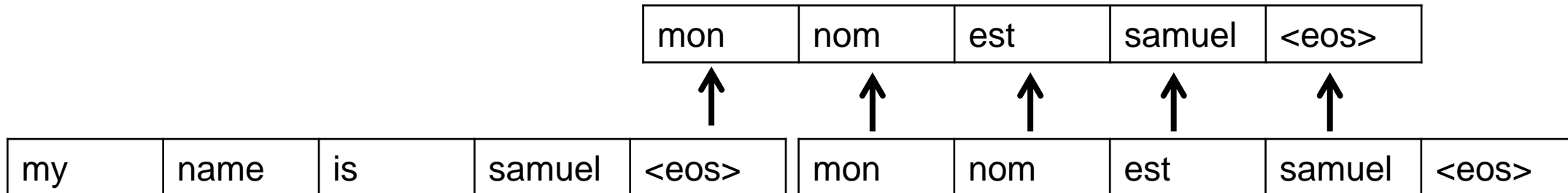
Sequence-to-sequence batch loss

Instead, we pass in all the English (or whichever) tokens, and then do teacher forcing loss on the French (or whatever) tokens only.

So for:



The loss would look like:



Preliminaries



Similar preliminaries:

1. Download translation dataset:
 - <https://download.pytorch.org/tutorial/data.zip>
2. Load 2 vector models-one for English, one for French
 - Add special tokens to both: <pad>, <unk>, <sos>, <eos>
3. Preprocess translation dataset and map to vector model tokens
4. Create dataset & dataloader
5. Install PyTorch Lightning

	english	french
0	Go.	Va !
1	Run!	Cours !
2	Run!	Courez !
3	Wow!	Ça alors !
4	Fire!	Au feu !
5	Help!	À l'aide !
6	Jump.	Saute.
7	Stop!	Ça suffit !
8	Stop!	Stop !
9	Stop!	Arrête-toi !



Preprocessed dataset

```
1 eng_fra_df.head(10)
```

	english	french	english_tokens	french_tokens	english_token_ids	french_token_ids
0	Go.	Va !	[<sos>, go, ., <eos>]	[<sos>, va, !, <eos>]	[400002, 242, 2, 400003]	[155564, 158, 155562, 155565]
1	Run!	Cours !	[<sos>, run, !, <eos>]	[<sos>, cours, !, <eos>]	[400002, 307, 805, 400003]	[155564, 239, 155562, 155565]
2	Run!	Courez !	[<sos>, run, !, <eos>]	[<sos>, courez, !, <eos>]	[400002, 307, 805, 400003]	[155564, 38881, 155562, 155565]
3	Wow!	Ça alors !	[<sos>, wow, !, <eos>]	[<sos>, ça, alors, !, <eos>]	[400002, 14397, 805, 400003]	[155564, 110, 140, 155562, 155565]
4	Fire!	Au feu !	[<sos>, fire, !, <eos>]	[<sos>, au, feu, !, <eos>]	[400002, 484, 805, 400003]	[155564, 22, 1092, 155562, 155565]
5	Help!	À l'aide !	[<sos>, help, !, <eos>]	[<sos>, à, l'aide, !, <eos>]	[400002, 275, 805, 400003]	[155564, 7, 16685, 155562, 155565]
6	Jump.	Saute.	[<sos>, jump, ., <eos>]	[<sos>, saute, ., <eos>]	[400002, 3106, 2, 400003]	[155564, 11775, 155562, 155565]
7	Stop!	Ça suffit !	[<sos>, stop, !, <eos>]	[<sos>, ça, suffit, !, <eos>]	[400002, 837, 805, 400003]	[155564, 110, 1292, 155562, 155565]
8	Stop!	Stop !	[<sos>, stop, !, <eos>]	[<sos>, stop, !, <eos>]	[400002, 837, 805, 400003]	[155564, 6517, 155562, 155565]
9	Stop!	Arrête-toi !	[<sos>, stop, !, <eos>]	[<sos>, arrête-toi, !, <eos>]	[400002, 837, 805, 400003]	[155564, 155562, 155562, 155565]



LSTM seq-to-seq model: `__init__()`

```
class EnglishFrenchSeqToSeqModel(pl.LightningModule):
    def __init__(self,
                 english_word_vectors:np.ndarray,
                 french_word_vectors:np.ndarray,
                 english_vocab_size:int,
                 french_vocab_size:int,
                 learning_rate:float,
                 english_padding_id:int,
                 french_padding_id:int,
                 french_eos_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,
                 loss_print_interval=100,
                 **kwargs):
        super().__init__(**kwargs)

        self.english_embeddings = torch.nn.Embedding.from_pretrained(torch.tensor(english_word_vectors),
                                                                    freeze=True)
        self.french_embeddings = torch.nn.Embedding.from_pretrained(torch.tensor(french_word_vectors),
                                                                    freeze=True)
        self.lstm = torch.nn.LSTM(input_size = english_word_vectors.shape[1], # The LSTM will be taking in word vectors
                                  hidden_size = lstm_hidden_size,
                                  num_layers=lstm_layers,
                                  bidirectional=False, # We can't count on being able to proceed both backward and forward
                                  dropout=dropout_prob,
                                  batch_first=True # This is important. Set to False by default for some reason.
                                  )

        # Output layer has to produce one logit per potential word, so the output size is vocab_size
        self.output_layer = torch.nn.Linear(lstm_hidden_size, french_vocab_size)
        self.lstm_layers = lstm_layers
        self.learning_rate = learning_rate
        self.english_padding_id = english_padding_id
        self.french_padding_id = french_padding_id
        self.french_eos_id = french_eos_id
```



LSTM seq-to-seq model: forward()

```
def forward(self,
            english_input_ids:torch.Tensor, #(batch size x max sequence length),
            french_input_ids:torch.Tensor
            ):

    # First we need to pass the English tokens into the model to generate a final hidden vector and cell state
    english_embeds = self.english_embeddings(english_input_ids) #(batch size x max english sequence length x embedding size)
    english_padding_mask = (english_input_ids != self.english_padding_id).int() #(batch size x max english sequence length)
    english_input_lengths = english_padding_mask.sum(dim=1).detach().cpu() #(batch size)
    english_packed_embeddings = pack_padded_sequence(english_embeds, english_input_lengths, batch_first=True, enforce_sorted=False)
    english_packed_output, final_english_hidden, final_english_state = self.lstm.forward(english_packed_embeddings)

    # We don't need to bother unpacking the English outputs since we're not using them
    # english_hiddens, _ = pad_packed_sequence(english_packed_output, batch_first=True, padding_value=0.0, total_length=english_input_ids.shape[1])

    # Then we start with that final hidden and final state, and do teacher-forcing using the french sequence vs. itself
    french_embeds = self.french_embeddings(french_input_ids) #(batch size x max french sequence length x embedding size)
    french_padding_mask = (french_input_ids != self.french_padding_id).int() #(batch size x max french sequence length)
    french_input_lengths = french_padding_mask.sum(dim=1).detach().cpu() #(batch size)
    french_packed_embeddings = pack_padded_sequence(french_embeds, french_input_lengths, batch_first=True, enforce_sorted=False)
    french_packed_output, (final_french_hidden, final_french_state) = self.lstm.forward(french_packed_embeddings, (final_english_hidden, final_english_state))
    french_hiddens, _ = pad_packed_sequence(french_packed_output, batch_first=True, padding_value=0.0, total_length=french_input_ids.shape[1])

    french_output_logits = self.output_layer(french_hiddens) #(batch size x max french sequence length x vocab size)

    losses = torch.nn.functional.cross_entropy(french_output_logits[:, :-1].transpose(1,2), french_input_ids[:, 1:], reduction='none')

    # Then the final thing we need to do is zero out the losses whenever the target token is a padding token
    padded_losses = losses * french_padding_mask[:, 1:] # (batch size x max french sequence length)
    loss = padded_losses.mean() #(1)

    return {'loss':loss}
```



LSTM seq-to-seq model: generate()

```
def generate(self,
             english_input_ids:torch.Tensor, # shape (1,sequence length) or (sequence length)
             max_output_length:int, # How many tokens to generate past the input sequence
             temperature:float=0.5, # How loosely to sample from the output distribution
             ):

    # If the input shape is (1, sequence length), make it (sequence length)
    if english_input_ids.ndim == 2: english_input_ids = english_input_ids.squeeze(0)

    # Remove padding tokens if they are present
    english_padding_mask = (english_input_ids != self.english_padding_id).int() #(batch size x max sequence length)
    english_input_length = english_padding_mask.sum().detach().cpu() #(batch size)
    english_input_ids = english_input_ids[0:english_input_length]
    english_inputs_embeds = self.english_embeddings(english_input_ids) #(sequence length x embedding size)

    # First we run the given sequence through the LSTM
    # Because we aren't using a batch of variable-length sequences, we don't have to bother with a packed padded sequence like above
    english_hidden, (final_english_hidden, final_english_state) = self.lstm.forward(english_inputs_embeds) # (sequence length x lstm hidden size),
    #( lstm layers x lstm hidden size), (lstm layers x lstm hidden size)

    output_tokens = []
    # Then, for the rest of the desired output length, we generate one token at a time, conditioned on the previous generated token
    last_hidden, last_state = final_english_hidden, final_english_state
    last_logits = self.output_layer(last_hidden[-1])
    last_token_id = self.sample_token_id_from_logits(last_logits, temperature)
    output_tokens.append(last_token_id)
    for i in range(english_input_length, max_output_length):
        last_embeds = self.french_embeddings(last_token_id).unsqueeze(0)
        last_output, (last_hidden, last_state) = self.lstm.forward(last_embeds, (last_hidden, last_state))
        last_logits = self.output_layer(last_hidden[-1])
        last_token_id = self.sample_token_id_from_logits(last_logits, temperature)
        output_tokens.append(last_token_id)
        # Uncomment this if you want the model to stop on <eos>
        # if last_token_id == self.french_eos_id:
        #     break

    output_ids = torch.stack(output_tokens)

    return {'french_output_ids':output_ids}
```



Train model

```
1 seqtoseq_model = EnglishFrenchSeqToSeqModel(  
2     english_word_vectors=english_vector_model.vectors,  
3     english_vocab_size = english_vector_model.vectors.shape[0],  
4     french_word_vectors=french_vector_model.vectors,  
5     french_vocab_size = french_vector_model.vectors.shape[0],  
6     learning_rate = 0.001, #I'll typically start with something like 1e-3 for LSTMs  
7     english_padding_id = english_vector_model.key_to_index['<pad>'],  
8     french_padding_id = french_vector_model.key_to_index['<pad>'],  
9     french_eos_id = french_vector_model.key_to_index['<eos>'],  
10    lstm_hidden_size=100,  
11    lstm_layers=2,  
12    dropout_prob=0.1)  
13  
14 from pytorch_lightning import Trainer  
15 from pytorch_lightning.callbacks.progress import TQDMProgressBar  
16  
17 trainer = Trainer(  
18     accelerator="auto",  
19     devices=1 if torch.cuda.is_available() else None,  
20     max_epochs=5,  
21     callbacks=[TQDMProgressBar(refresh_rate=20)],  
22     val_check_interval = 0.1,  
23 )  
24 trainer.fit(model=seqtoseq_model,  
25             train_dataloaders=train_dataloader,  
26             val_dataloaders=val_dataloader)
```

```
Mean training loss (steps -100-0): nan  
Mean training loss (steps 0-100): 4.271  
Mean training loss (steps 100-200): 3.373  
Epoch 0 step 250 validation loss: tensor(3.2704, device='cuda:0')  
Mean training loss (steps 200-300): 3.250  
Mean training loss (steps 300-400): 3.109  
Epoch 0 step 500 validation loss: tensor(2.9881, device='cuda:0')  
Mean training loss (steps 400-500): 3.076  
Mean training loss (steps 500-600): 2.968  
Mean training loss (steps 600-700): 2.906  
Epoch 0 step 750 validation loss: tensor(2.8339, device='cuda:0')  
Mean training loss (steps 700-800): 2.901  
Mean training loss (steps 800-900): 2.844  
Epoch 0 step 1000 validation loss: tensor(2.7491, device='cuda:0')
```

```
Epoch 4 step 11750 validation loss: tensor(1.7015, device='cuda:0')  
Mean training loss (steps 1700-1800): 1.540  
Mean training loss (steps 1800-1900): 1.558  
Epoch 4 step 12000 validation loss: tensor(1.6946, device='cuda:0')  
Mean training loss (steps 1900-2000): 1.612  
Mean training loss (steps 2000-2100): 1.505  
Mean training loss (steps 2100-2200): 1.572  
Epoch 4 step 12250 validation loss: tensor(1.6841, device='cuda:0')  
Mean training loss (steps 2200-2300): 1.564  
Mean training loss (steps 2300-2400): 1.562  
Epoch 4 step 12500 validation loss: tensor(1.6682, device='cuda:0')  
INFO:pytorch_lightning.utilities.rank_zero:`Trainer.fit` stopped: `max_epochs=
```




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

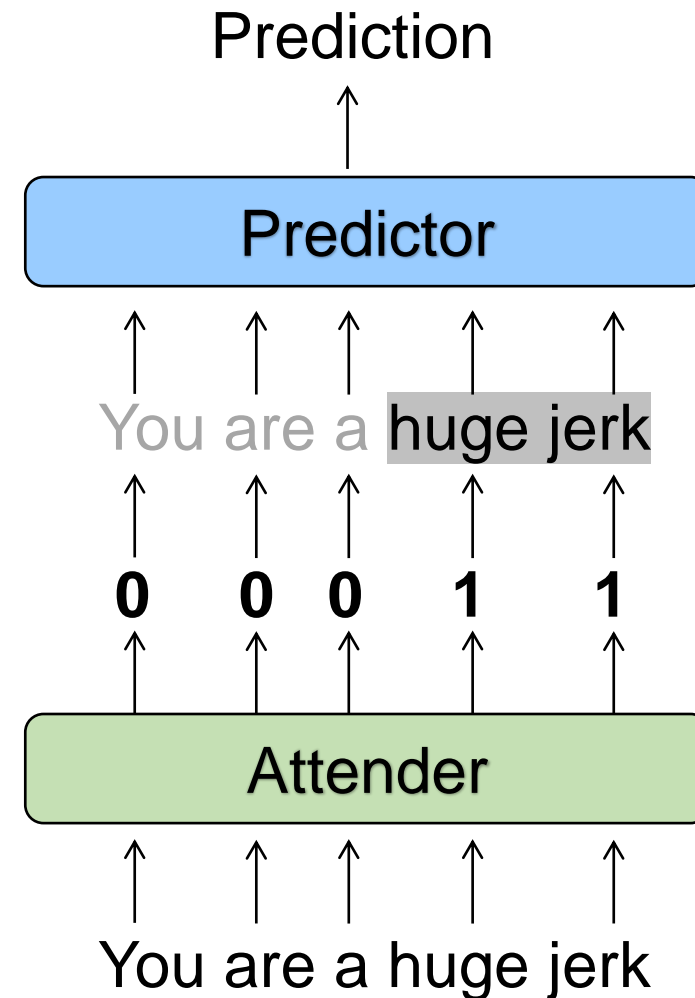
Solution: **Attention**

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.





Attention classification model

```
1 class AttentionClassifier(pl.LightningModule):
2     def __init__(self,
3         word_vectors:np.ndarray,
4         num_classes:int,
5         learning_rate:float,
6         padding_id:int,
7         lstm_hidden_size:int=100, # how big the inner vectors of the LSTM will be,
8         lstm_layers:int =2, # how many layers the LSTM will have
9         dropout_prob:float=0.1,
10        sparsity_loss_weight:float= 0.15,
11        **kwargs):
12     super().__init__( **kwargs)
13
14     # We'll use the same PyTorch Embedding layer as before
15     self.word_embeddings = torch.nn.Embedding.from_pretrained(torch.tensor(word_vectors),
16                                                                freeze=True)
17
18
19     self.attender = torch.nn.LSTM(input_size = word_vectors.shape[1],
20                                  hidden_size = lstm_hidden_size,
21                                  num_layers=lstm_layers,
22                                  bidirectional=True,
23                                  dropout=dropout_prob,
24                                  batch_first=True)
25     self.attender_output_layer = torch.nn.Linear(2*lstm_hidden_size, 1)
26
27     self.predictor = torch.nn.LSTM(input_size = word_vectors.shape[1],
28                                    hidden_size = lstm_hidden_size,
29                                    num_layers=lstm_layers,
30                                    bidirectional=True,
31                                    dropout=dropout_prob,
32                                    batch_first=True)
33     self.predictor_output_layer = torch.nn.Linear(2*lstm_hidden_size, num_classes)
34
35
36     # Output layer input size has to be doubled because the LSTM is bidirectional
37     self.lstm_layers = lstm_layers
38     self.learning_rate = learning_rate
39     self.padding_id = padding_id # we'll need this later
40     self.sparsity_loss_weight = sparsity_loss_weight
41     self.train_accuracy = Accuracy(task='multiclass', num_classes=num_classes)
42     self.val_accuracy = Accuracy(task='multiclass', num_classes=num_classes)
```



Attention classification model

```
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

```
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_data loaders=train_data loader,  
3 | | | | | val_data loaders=dev_data loader)
```

/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

```
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])
```

```
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
.
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
```




Concluding thoughts

Sequence-to-sequence models

- Main application: translation

Attention

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

Model saving/loading