



Deep Learning School

Физтех-Школа Прикладной математики и информатики (ФПМИ)
МФТИ

Some parts of the notebook are almost the exact copy of
https://github.com/yandexdataschool/nlp_course

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

Attention layer can take in the previous hidden state of the decoder s_{t-1} , and all of the stacked forward and backward hidden states H from the encoder. The layer will output an attention vector a_t , that is the length of the source sentence, each element is between 0 and 1 and the entire vector sums to 1.

Intuitively, this layer takes what we have decoded so far s_{t-1} , and all of what we have encoded H , to produce a vector a_t , that represents which words in the source sentence we should pay the most attention to in order to correctly predict the next word to decode \hat{y}_{t+1} . The decoder input word that has been embedded y_t .

You can use any type of the attention scores between previous hidden state of the encoder s_{t-1} and hidden state of the decoder $h \in H$, you prefer. We have met at least three of them:

$$\text{score}(\mathbf{h}, \mathbf{s}_{t-1}) = \begin{cases} \mathbf{h}^\top \mathbf{s}_{t-1} & \text{dot} \\ \mathbf{h}^\top \mathbf{W}_a \mathbf{s}_{t-1} & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}; \mathbf{s}_{t-1}]) & \text{concat} \end{cases}$$

We wil use "concat attention":

First, we calculate the *energy* between the previous decoder hidden state \mathbf{s}_{t-1} and the encoder hidden states \mathbf{H} . As our encoder hidden states \mathbf{H} are a sequence of T tensors, and our previous decoder hidden state \mathbf{s}_{t-1} is a single tensor, the first thing we do is repeat the previous decoder hidden state T times. \Rightarrow

We have:

$$\mathbf{H} = [\mathbf{h}_0, \dots, \mathbf{h}_{T-1}]$$

$$[\mathbf{s}_{t-1}, \dots, \mathbf{s}_{t-1}]$$

The encoder hidden dim and the decoder hidden dim should be equal: **dec hid dim = enc hid dim**.

We then calculate the energy, E_t , between them by concatenating them together:

$$[[\mathbf{h}_0, \mathbf{s}_{t-1}], \dots, [\mathbf{h}_{T-1}, \mathbf{s}_{t-1}]]$$

And passing them through a linear layer (attn = \mathbf{W}_a) and a \tanh activation function:

$$E_t = \tanh(\text{attn}(\mathbf{H}, \mathbf{s}_{t-1}))$$

This can be thought of as calculating how well each encoder hidden state "matches" the previous decoder hidden state.

We currently have a **[enc hid dim, src sent len]** tensor for each example in the batch. We want

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu on should be over the length by a **[1, enc hid dim]** tensor, \mathbf{v}

$$\hat{a}_t = \mathbf{v} E_t$$

We can think of this as calculating a weighted sum of the "match" over all enc_hid_dem elements for each encoder hidden state, where the weights are learned (as we learn the parameters of \mathbf{v}).

Finally, we ensure the attention vector fits the constraints of having all elements between 0 and 1 and the vector summing to 1 by passing it through a softmax layer.

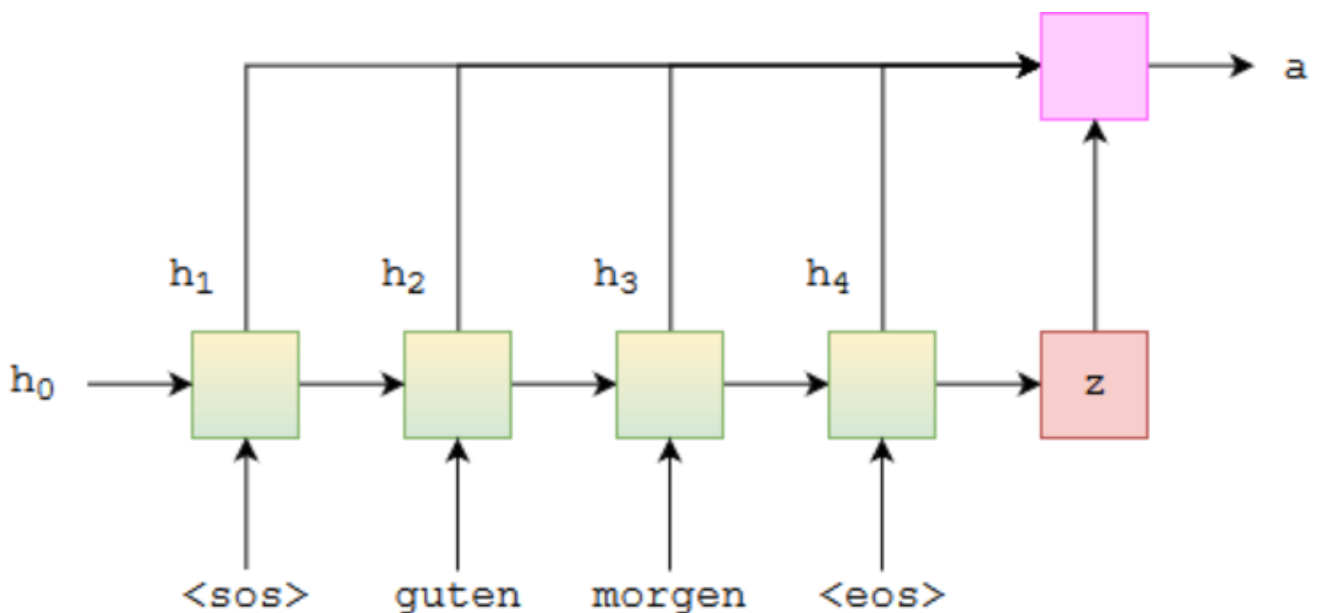
$$a_t = \text{softmax}(\hat{a}_t)$$

▼ Temperature SoftMax

$$\text{softmax}(x)_i = \frac{e^{\frac{y_i}{T}}}{\sum_j^N e^{\frac{y_j}{T}}}$$

This gives us the attention over the source sentence!

Graphically, this looks something like below. $z = s_{t-1}$. The green/yellow blocks represent the hidden states from both the forward and backward RNNs, and the attention computation is all done within the pink block.



Neural Machine Translation

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ☐ h convergence
plots/metrics and your thoughts. Just like you would approach a real problem.

```
1 ! wget https://drive.google.com/uc?id=1NWYqJgeG_4883LINdEjKUr6nLQPY6Yb_-O data.txt
2
3 # Thanks to YSDA NLP course team for the data
4 # (who thanks tilda and deephack teams for the data in their turn)
```

```
--2021-11-07 10:01:40-- https://drive.google.com/uc?id=1NWYqJgeG_4883LINdEjKUr6nLQP
Resolving drive.google.com (drive.google.com)... 209.85.147.138, 209.85.147.139, 209
Connecting to drive.google.com (drive.google.com)|209.85.147.138|:443... connected.
HTTP request sent, awaiting response... 302 Moved Temporarily
Location: https://doc-14-00-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717
Warning: wildcards not supported in HTTP.
--2021-11-07 10:01:40-- https://doc-14-00-docs.googleusercontent.com/docs/securesc/
Resolving doc-14-00-docs.googleusercontent.com (doc-14-00-docs.googleusercontent.com)
Connecting to doc-14-00-docs.googleusercontent.com (doc-14-00-docs.googleusercontent
HTTP request sent, awaiting response... 200 OK
Length: 12905334 (12M) [text/plain]
```

Saving to: 'data.txt'

data.txt 100%[=====>] 12.31M --.-KB/s in 0.08s

2021-11-07 10:01:41 (163 MB/s) - 'data.txt' saved [12905334/12905334]



```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4
5 import torchtext
6 from torchtext.legacy.data import Field, BucketIterator
7
8 import spacy
9
10 import random
11 import math
12 import time
13 import numpy as np
14
15 import matplotlib
16 matplotlib.rcParams.update({'figure.figsize': (16, 12), 'font.size': 14})
17 import matplotlib.pyplot as plt
18 %matplotlib inline
19 from IPython.display import clear_output
20
21 from nltk.tokenize import WordPunctTokenizer
```

We'll set the random seeds for deterministic results.

```
1 SEED = 1234
```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```
5 torch.manual_seed(SEED)
6 torch.cuda.manual_seed(SEED)
7 torch.backends.cudnn.deterministic = True
```

▼ Preparing Data

Here comes the preprocessing

```
1 tokenizer_W = WordPunctTokenizer()
2
3 def tokenize_ru(x, tokenizer=tokenizer_W):
4     return tokenizer.tokenize(x.lower())[:-1]
5
6 def tokenize_en(x, tokenizer=tokenizer_W):
7     return tokenizer.tokenize(x.lower())
```

```

1 tokenize_ru("Привет, как дела")

    ['дела', 'как', ',', 'привет']

1 SRC = Field(tokenize=tokenize_ru,
2             init_token = '<sos>',
3             eos_token = '<eos>',
4             lower = True)
5
6 TRG = Field(tokenize=tokenize_en,
7             init_token = '<sos>',
8             eos_token = '<eos>',
9             lower = True)
10
11
12 dataset = torchtext.legacy.data.TabularDataset(
13     path='data.txt',
14     format='tsv',
15     fields=[('trg', TRG), ('src', SRC)]
16 )

1 print(len(dataset.examples))
2 print(dataset.examples[0].src)
3 print(dataset.examples[0].trg)

50000
['.', 'собора', 'троицкого', '-', 'свято', 'от', 'ходьбы', 'минутах', '3', 'в', ',',
['cordelia', 'hotel', 'is', 'situated', 'in', 'tbilisi', ',', 'a', '3', '-', 'minute

1 train_data, valid_data, test_data = dataset.split(split_ratio=[0.8, 0.15, 0.05])
2
3 print(f"Number of training examples: {len(train_data.examples)}")
4 print(f"Number of validation examples: {len(valid_data.examples)}")
5 print(f"Number of testing examples: {len(test_data.examples)}")

Number of training examples: 40000
Number of validation examples: 2500
Number of testing examples: 7500

1 SRC.build_vocab(train_data, min_freq = 2)
2 TRG.build_vocab(train_data, min_freq = 2)

1 print(f"Unique tokens in source (ru) vocabulary: {len(SRC.vocab)}")
2 print(f"Unique tokens in target (en) vocabulary: {len(TRG.vocab)}")

Unique tokens in source (ru) vocabulary: 14129
Unique tokens in target (en) vocabulary: 10104

```

And here is example from train dataset:

```
1 print(vars(train_data.examples[9]))

{'trg': ['other', 'facilities', 'offered', 'at', 'the', 'property', 'include', 'groc
```

When we get a batch of examples using an iterator we need to make sure that all of the source sentences are padded to the same length, the same with the target sentences. Luckily, TorchText iterators handle this for us!

We use a `BucketIterator` instead of the standard `Iterator` as it creates batches in such a way that it minimizes the amount of padding in both the source and target sentences.

```
1 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

```
1 !nvidia-smi
```

```
Sun Nov  7 10:02:00 2021
```

NVIDIA-SMI 495.44				Driver Version: 460.32.03				CUDA Version: 11.2	

GPU	Name		Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC	
Fan	Temp	Perf	Pwr:Usage/Cap		Memory-Usage	GPU-Util	Compute	M.	
								MIG M.	
=====									
0	Tesla K80		Off	00000000:00:04.0	Off			0	
N/A	69C	P8	31W / 149W		3MiB / 11441MiB	0%	Default	N/A	

Processes:									
GPU	GI	CI	PID	Type	Process name	GPU Memory Usage			
	ID	ID							
=====									

Undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu

```
1 def _len_sort_key(x):
2     return len(x.src)
3
4 BATCH_SIZE = 128
5
6 train_iterator, valid_iterator, test_iterator = BucketIterator.splits(
7     (train_data, valid_data, test_data),
8     batch_size = BATCH_SIZE,
9     device = device,
10    sort_key=_len_sort_key
11 )
```

▼ Let's use modules.py

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.

```
1 %cd /content/drive/MyDrive/Colab Notebooks/Attention
```

/content/drive/MyDrive/Colab Notebooks/Attention

```
1 %ls
```

```
best-val-model.pt      modules.py
'[homework]NeuralMachineTranslation.ipynb'  __pycache__/
```

```
1 %cd ../drive/MyDrive/Colab Notebooks/Attention/modules.py
```

[Errno 2] No such file or directory: '../drive/MyDrive/Colab Notebooks/Attention/modu
/content/drive/MyDrive/Colab Notebooks/Attention

▼ Encoder

For a multi-layer RNN, the input sentence, X , goes into the first (bottom) layer of the RNN and hidden states, $H = \{h_1, h_2, \dots, h_T\}$, output by this layer are used as inputs to the RNN in the layer above. Thus, representing each layer with a superscript, the hidden states in the first layer are given by:

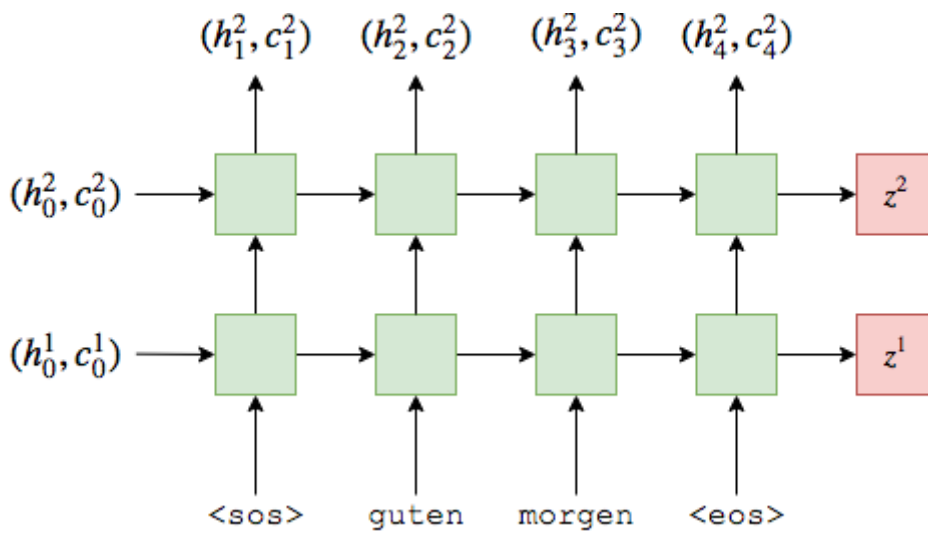
$$h_t^1 = \text{EncoderRNN}^1(x_t, h_{t-1}^1)$$

The hidden states in the second layer are given by:

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

Extending our multi-layer equations to LSTMs, we get:

$$\begin{aligned}(h_t^1, c_t^1) &= \text{EncoderLSTM}^1(x_t, (h_{t-1}^1, c_{t-1}^1)) \\ (h_t^2, c_t^2) &= \text{EncoderLSTM}^2(h_t^1, (h_{t-1}^2, c_{t-1}^2))\end{aligned}$$



```

1 # you can paste code of encoder from modules.py
2 # the encoder can be like seminar encoder but you have to return outputs
3 # and if you use bidirectional you won't make the same operation like with hidden
4 # because outputs = [src sent len, batch size, hid dim * n directions]
5 class Encoder(nn.Module):
6     def __init__(self, input_dim, emb_dim, hid_dim, n_layers, dropout, bidirectional):
7         super().__init__()
8
9         self.input_dim = input_dim           # source vocab size
10        self.emb_dim = emb_dim                 # vector length for each token
11        self.hid_dim = hid_dim                 # hidden dim of the RNN
12        self.n_layers = n_layers               # num LSTM layers
13        self.dropout = dropout                 # prob of zeroing out units
14        self.bidirectional = bidirectional     # bool, if the RNN will be bidirectional
15
16        self.embedding = nn.Embedding(input_dim, emb_dim)
17
18        self.rnn = nn.LSTM(emb_dim, hid_dim, num_layers=n_layers, dropout=dropout, bidi
19
20
21
22    def forward(self, src):
23
24        #src = [src sent len, batch size]
25
26        # Compute an embedding from the src data and apply dropout to it
27        embedded = self.dropout(self.embedding(src))
28
29        #embedded = [src sent len, batch size, emb dim]
30
31        # Compute the RNN output values of the encoder RNN.
32        # outputs, hidden and cell should be initialized here. Refer to nn.LSTM docs ;)
33        # https://pytorch.org/docs/stable/generated/torch.nn.LSTM.html
34
35        outputs, (hidden, cell) = self.rnn(embedded)
36        #print("outputs.shape\n\t", outputs.shape)
37
38        #outputs = [src sent len, batch size, hid dim * n directions]

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕


```

39
40     #hidden = [n layers * n directions, batch size, hid dim] | for n directions =
41     #         = [n layers, batch size, hid dim]              | for n directions =
42     #         = [n layers, batch size, hid dim * n directions]
43     #cell = [n layers * n directions, batch size, hid dim] | for n directions =
44     #       = [n layers, batch size, hid dim]              | for n directions =
45     #       = [n layers, batch size, hid dim * n directions]
46
47
48     #outputs are always from the top hidden layer
49     if self.bidirectional:
50         #print("hidden.shape before adjustment\n\t", hidden.shape)
51         hidden = hidden.reshape(self.n_layers, 2, -1, self.hid_dim)
52         hidden = hidden.transpose(1, 2).reshape(self.n_layers, -1, 2 * self.hid_dim)
53         #print("hidden.shape after adjustment\n\t", hidden.shape)
54
55         #print("cell.shape before adjustment\n\t", cell.shape)
56         cell = cell.reshape(self.n_layers, 2, -1, self.hid_dim)
57         cell = cell.transpose(1, 2).reshape(self.n_layers, -1, 2 * self.hid_dim)
58         #print("cell.shape after adjustment\n\t", cell.shape)
59
60
61     # in both cases bidirectional=True/False we get the followong shapes
62     #   hidden = [n layers, batch size, hid dim * n directions]
63     #   cell = [n layers, batch size, hid dim * n directions]
64     return outputs, hidden, cell

```

▼ Attention

$$\text{score}(\mathbf{h}, \mathbf{s}_{t-1}) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}; \mathbf{s}_{t-1}]) - \text{concat attention}$$

```

1 # you can paste code of attention from modules.py
2
5 #     return e_x / torch.sum(e_x, dim=0)
6
7
8 # use your temperature
9 def softmax(x, temperature):
10     e_x = torch.exp(x / temperature)
11     return e_x / torch.sum(e_x, dim=0)
12
13
14
15
16 class Attention(nn.Module):
17     def __init__(self, enc_hid_dim, dec_hid_dim, softmax_temp=1):
18         super().__init__()
19
20         self.enc_hid_dim = enc_hid_dim
21         self.dec_hid_dim = dec_hid_dim

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

22
23     self.attn = nn.Linear(enc_hid_dim + dec_hid_dim, enc_hid_dim)
24     self.v = nn.Linear(enc_hid_dim, 1)
25
26     # defaults to 1 (the usual softmax w/o temperature)
27     self.softmax_temp = softmax_temp
28
29     def forward(self, hidden, encoder_outputs):
30
31         # encoder_outputs = [src sent len, batch size, enc_hid_dim]
32
33         # only take the last layer of the decoder's hidden units
34         # hidden = [n_layers, batch size, dec_hid_dim]
35         last_hidden = hidden[-1, :, :].unsqueeze(0)
36         # last_hidden = [1, batch size, dec_hid_dim]
37
38         # repeat hidden and concatenate it with encoder_outputs
39         hiddens = last_hidden.repeat(encoder_outputs.shape[0], 1, 1)
40         # hiddens = [src sent len, batch size, dec_hid_dim]
41
42         #print("encoder_outputs.shape:", encoder_outputs.shape)
43         #print("hiddens.shape:", hiddens.shape)
44
45         concat_h_s = torch.cat([hiddens, encoder_outputs], dim=2)
46         #print("concat_h_s.shape:", concat_h_s.shape)
47         #print("(self.enc_hid_dim + self.dec_hid_dim).shape:", self.enc_hid_dim + self.
48         # concat_h_s = [src sent len, batch size, enc_hid_dim + dec_hid_dim]
49
50         # calculate energy: E
51         E = torch.tanh(self.attn(concat_h_s))
52         # E = [src sent len, batch size, enc_hid_dim] (see self.attn)
53
54         # get attention (not normalized to probabilities yet)
55         a = self.v(E)
56         # a = [src sent len, batch size, 1]

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

59     return softmax(a, temperature=self.softmax_temp)

```

```

1 # src sent len = 3
2 # batch_size = 8
3 # enc_hid_dim = 2
4 # dec_hid_dim = 3
5
6 d = torch.ones(3, 8, 2) * 2    # encoder_outputs = [src sent len, batch size, enc_hid_
7 c = torch.ones(10, 8, 3)      # hidden = [n_layers, batch size, dec_hid_dim]
8
9 print(c[-1, :, :].shape)
10
11 c[-1, :, :].unsqueeze(0).shape

```

```

    torch.Size([8, 3])
    torch.Size([1, 8, 3])

```

▼ Decoder with Attention

To make it really work you should also change the `Decoder` class from the classwork in order to make it to use `Attention`. You may just copy-paste `Decoder` class and add several lines of code to it.

The decoder contains the attention layer `attention`, which takes the previous hidden state s_{t-1} , all of the encoder hidden states H , and returns the attention vector a_t .

We then use this attention vector to create a weighted source vector, w_t , denoted by `weighted`, which is a weighted sum of the encoder hidden states, H , using a_t as the weights.

$$w_t = a_t H$$

The input word that has been embedded y_t , the weighted source vector w_t , and the previous decoder hidden state s_{t-1} , are then all passed into the decoder RNN, with y_t and w_t being concatenated together.

$$s_t = \text{DecoderGRU}([y_t, w_t], s_{t-1})$$

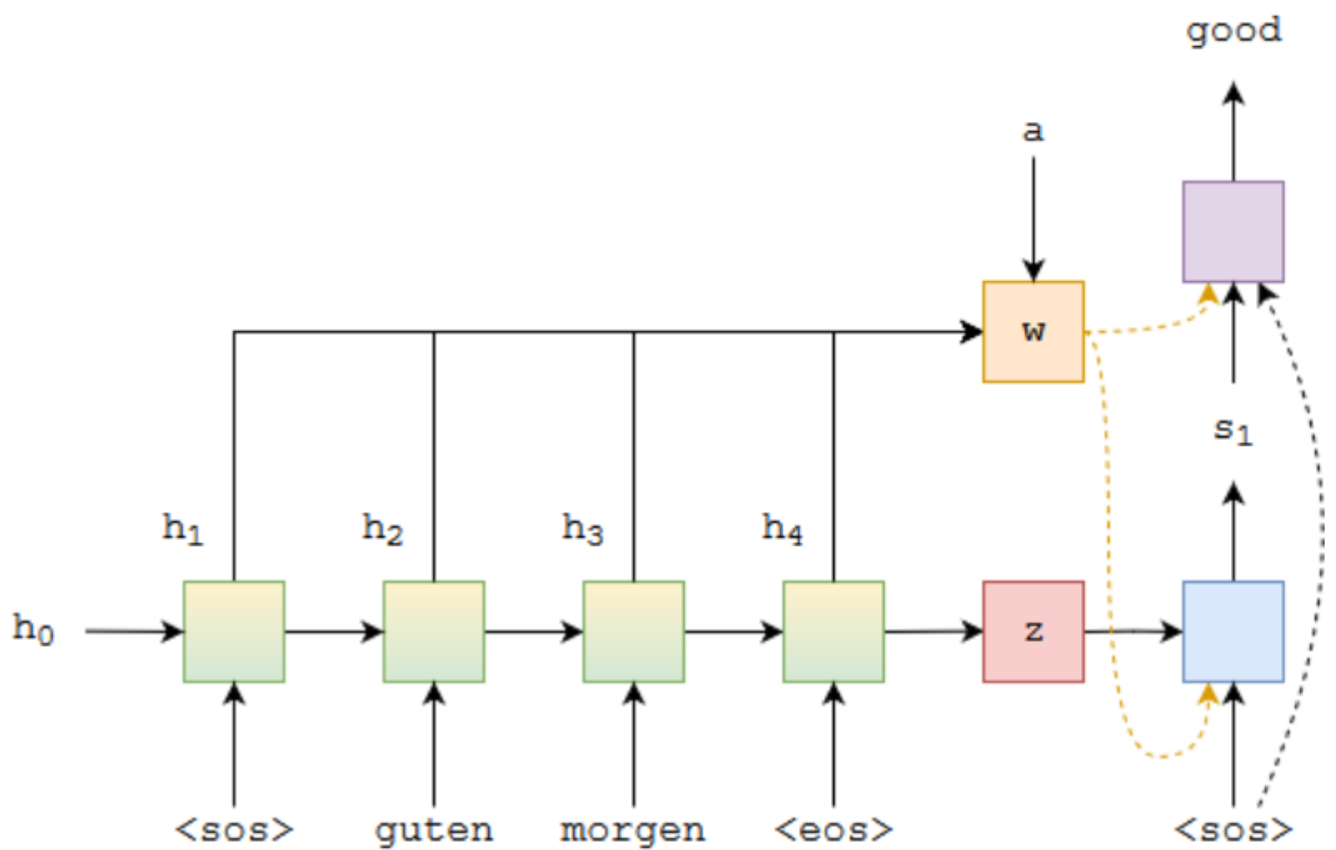
We then pass y_t , w_t and s_t through the linear layer, f , to make a prediction of the next word in the target sentence, \hat{y}_{t+1} . This is done by concatenating them all together.

$$\hat{y}_{t+1} = f(y_t, w_t, s_t)$$

The image below shows decoding the **first** word in an example translation.

The green/yellow blocks show the forward/backward encoder RNNs which output H , the red block is $z = s_{t-1} = s_0$, the blue block shows the decoder RNN which outputs $s_t = s_1$, the purple block shows the linear layer, f , which outputs \hat{y}_{t+1} and the orange block shows the calculation of the weighted sum over H by a_t and outputs w_t . Not shown is the calculation of a_t .

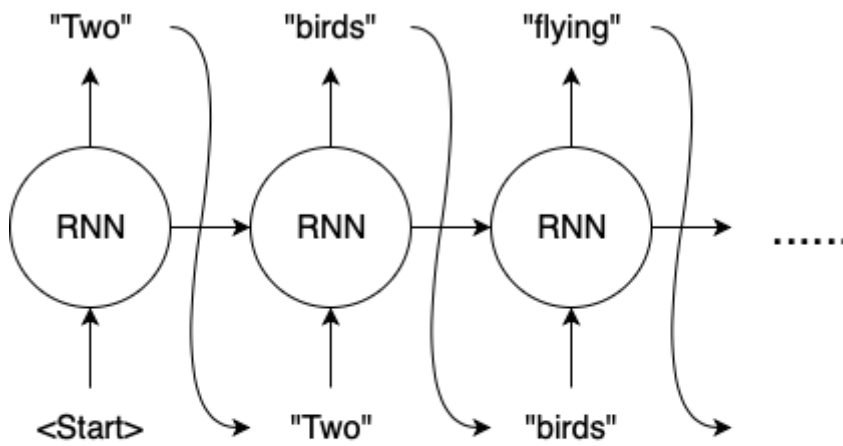
To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕



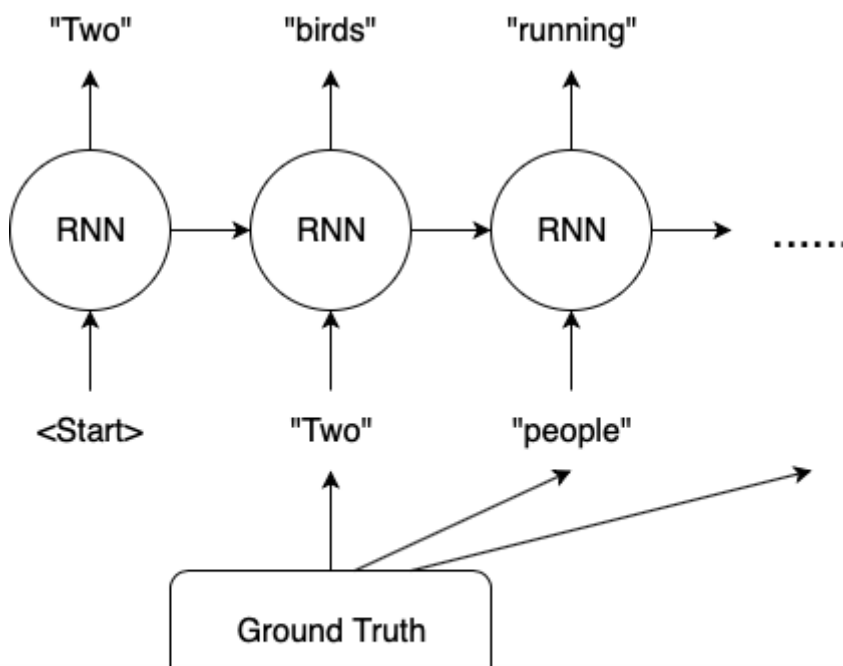
▼ Teacher forcing

Teacher forcing is a method for quickly and efficiently training recurrent neural network models that use the ground truth from a prior time step as input.

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕



Without Teacher Forcing



To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

With Teacher Forcing

When training/testing our model, we always know how many words are in our target sentence, so we stop generating words once we hit that many. During inference (i.e. real world usage) it is common to keep generating words until the model outputs an `<eos>` token or after a certain amount of words have been generated.

Once we have our predicted target sentence, $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_T\}$, we compare it against our actual target sentence, $Y = \{y_1, y_2, \dots, y_T\}$, to calculate our loss. We then use this loss to update all of the parameters in our model.

```
1 # you can paste code of decoder from modules.py
2
3
4 class DecoderWithAttention(nn.Module):
```

```

5     def __init__(self, output_dim, emb_dim, enc_hid_dim, dec_hid_dim,
6                   n_layers, dropout, attention):
7         super().__init__()
8
9         self.emb_dim = emb_dim
10        self.enc_hid_dim = enc_hid_dim
11        self.dec_hid_dim = dec_hid_dim
12        self.output_dim = output_dim          # vocab size in the target language (hence
13        self.attention = attention
14
15        self.embedding = nn.Embedding(output_dim, emb_dim)
16
17
18        # use GRU: https://pytorch.org/docs/stable/generated/torch.nn.GRU.html
19        self.rnn = nn.GRU(input_size=emb_dim + enc_hid_dim, # see blue box on the image
20                          hidden_size=dec_hid_dim,
21                          num_layers=n_layers, dropout=dropout)
22
23        # linear layer to get next word: f(y_t, w_t, s_t)
24        # see purple box on the image
25        self.out = nn.Linear(in_features=emb_dim + enc_hid_dim + dec_hid_dim,
26                              out_features=output_dim)
27
28        self.dropout = nn.Dropout(dropout)
29
30    def forward(self, input, hidden, encoder_outputs):
31        # input = [batch size]
32        # hidden = [n layers * n directions, batch size, dec_hid_dim]
33        # print(encoder_outputs.shape)
34
35        #n directions in the decoder will both always be 1, therefore:
36        #hidden = [n layers, batch size, dec_hid_dim]
37
38        # input: [batch size] -> [1, batch size]
39        input = input.unsqueeze(0) # because only one word, seq_len=1
40
41        # apply dropout to it
42        embedded = self.dropout(self.embedding(input))
43        #embedded = [1, batch size, emb dim]
44
45        # get weighted sum of (encoder_outputs = [src sent len, batch size, enc_hid_dim]
46        a = self.attention(hidden, encoder_outputs)      # a = [src sent len, batch size
47        weighted = torch.bmm(a.permute(1, 2, 0),         # [batch size, 1, src sent len]
48                             encoder_outputs.transpose(0, 1) # [batch size, src sent len, er
49                             ).transpose(0, 1) # w = [batch size, 1, enc_hid_dim] -> [1, batch
50
51        # print("embedded.shape:", embedded.shape)
52        # print("w.shape:", w.shape)
53
54        # concatenate weighted sum and embedded, break through the GRU
55        # embedded = [1, batch size, emb dim]
56        # weighted = [1, batch size, enc_hid_dim]
57        output, hidden = self.rnn(torch.cat([embedded, weighted], dim=2), hidden)
58        # output = [1, batch size, dec_hid_dim]
59        # hidden = [n layers, batch size, dec_hid_dim]

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu   by dropout to it

```

60     # so if we want to use hidden in prediction, we need to get the last layer:
61     # hidden[-1,:,:] = [batch size, dec_hid_dim]
62
63
64     # need the dimentions to agree for torch.cat() to work
65     #
66     # torch.cat([embedded, weighted, output], dim=2).squeeze(0)
67     # = [batch size, emb_dim + enc_hid_dim + dec_hid_dim] =
68     # torch.cat([embedded, weighted, hidden[-1,:,:].unsqueeze(0)], dim=2).squeeze(0)
69     # =
70     # torch.cat([embedded.squeeze(0), weighted.squeeze(0), hidden[-1,:,:]], dim=1)
71
72     # get predictions
73     #
74     # my initial version:
75     # prediction = self.out(torch.cat([embedded, weighted, output], dim=2).squeeze(0))
76     #
77     # original paper version:
78     prediction = self.out(torch.cat([embedded, # [1, batch size, emb dim]
79                                     weighted, # [1, batch size, enc_hid_dim]
80                                     # hidden[-1,:,:].unsqueeze(0) = [1, batch size, dec_hid_dim]
81                                     hidden[-1,:,:].unsqueeze(0)], dim=2).squeeze(0))
82     #prediction = [batch size, output dim]
83
84     # will be used as arguments again:
85     # (part of input (top1 or teacher forced correct), hidden)
86     return prediction, hidden

```

▼ Seq2Seq

Main idea:

- $w_t = a_t H$

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

Note: our decoder loop starts at 1, not 0. This means the 0th element of our outputs tensor remains all zeros. So our trg and outputs look something like:

$$\begin{aligned} \text{trg} &= [\langle \text{sos} \rangle, y_1, y_2, y_3, \langle \text{eos} \rangle] \\ \text{outputs} &= [0, \hat{y}_1, \hat{y}_2, \hat{y}_3, \langle \text{eos} \rangle] \end{aligned}$$

Later on when we calculate the loss, we cut off the first element of each tensor to get:

$$\begin{aligned} \text{trg} &= [y_1, y_2, y_3, \langle \text{eos} \rangle] \\ \text{outputs} &= [\hat{y}_1, \hat{y}_2, \hat{y}_3, \langle \text{eos} \rangle] \end{aligned}$$

```

1 # you can paste code of seq2seq from modules.py
2
3 class Seq2Seq(nn.Module):
4     def __init__(self, encoder, decoder, device):

```

```

5         super().__init__()
6
7         self.encoder = encoder
8         self.decoder = decoder
9         self.device = device
10
11         assert encoder.hid_dim * (1 + encoder.bidirectional) == decoder.dec_hid_dim, \
12             "Hidden dimensions of encoder and decoder must be equal!"
13         # assert encoder.n_layers == decoder.n_layers, \
14             # "Encoder and decoder must have equal number of layers!"
15
16     def forward(self, src, trg, teacher_forcing_ratio = 0.5):
17
18         # src = [src sent len, batch size]
19         # trg = [trg sent len, batch size]
20         # teacher_forcing_ratio is probability to use teacher forcing
21         # e.g. if teacher_forcing_ratio is 0.75 we use ground-truth inputs 75% of the t
22
23         # Again, now batch is the first dimension instead of zero
24         trg_len = trg.shape[0]          # trg = [trg sent len, batch size]
25         batch_size = trg.shape[1]
26         trg_vocab_size = self.decoder.output_dim
27
28         #tensor to store decoder outputs
29         outputs = torch.zeros(trg_len, batch_size, trg_vocab_size).to(self.device)
30
31         #last hidden state of the encoder is used as the initial hidden state of the de
32         encoder_outputs, hidden, cell = self.encoder(src)
33         #print('Shapes for encoder_outputs, hidden, cell')
34         #print(encoder_outputs.shape, hidden.shape, cell.shape, '\n')
35
36         #first input to the decoder is the <sos> tokens
37         input = trg[0,:]
38
39         for t in range(1, trg_len):
40
41             # receive output tensor (predictions) and new hidden
42             # output, hidden = self.decoder(input, hidden, encoder_outputs)
43
44             #place predictions in a tensor holding predictions for each token
45             outputs[t] = output
46             #decide if we are going to use teacher forcing or not
47             teacher_force = random.random() < teacher_forcing_ratio
48             #get the highest predicted token from our predictions
49             top1 = output.argmax(-1)
50             #if teacher forcing, use actual next token as next input
51             #if not, use predicted token
52             input = trg[t] if teacher_force else top1
53
54         return outputs
55

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

, and encoder_outputs


```

1 # For reloading
2 import modules
3 import imp
4 imp.reload(modules)
5
6 Encoder = modules.Encoder
7 Attention = modules.Attention
8 Decoder = modules.DecoderWithAttention
9 Seq2Seq = modules.Seq2Seq

1 INPUT_DIM = len(SRC.vocab)                # 30,000
2 OUTPUT_DIM = len(TRG.vocab)               # 30,000
3 ENC_EMB_DIM = 256                         # 620
4 DEC_EMB_DIM = 256                         # 620
5 DEC_HID_DIM = 512                         # 1000
6 N_LAYERS = 2                             # improvement (default = 1)
7 ENC_DROPOUT = 0.5
8 DEC_DROPOUT = 0.5
9 BIDIRECTIONAL = True                     # True
10 ENC_HID_DIM = DEC_HID_DIM // (1 + BIDIRECTIONAL) # 1000 fwd + 1000 bwd = 2000
11
12 enc = Encoder(INPUT_DIM, ENC_EMB_DIM, ENC_HID_DIM, N_LAYERS, ENC_DROPOUT, BIDIRECTIONAL)
13
14 # we do not give encoder's bidirectional argument to attention and decoder, but need to
15 att = Attention(enc_hid_dim=ENC_HID_DIM*(1+BIDIRECTIONAL), dec_hid_dim=DEC_HID_DIM, softmax=1)
16 dec = DecoderWithAttention(OUTPUT_DIM, DEC_EMB_DIM, ENC_HID_DIM*(1+BIDIRECTIONAL), DEC_HID_DIM,
17                             N_LAYERS, DEC_DROPOUT, BIDIRECTIONAL)
18 # dont forget to put the model to the right device
19 model = Seq2Seq(enc, dec, device).to(device)

```

```

1 def init_weights(m):
2     for name, param in m.named_parameters():

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

3 model.apply(init_weights)

Seq2Seq(
  (encoder): Encoder(
    (embedding): Embedding(14129, 256)
    (rnn): LSTM(256, 256, num_layers=2, dropout=0.5, bidirectional=True)
    (dropout): Dropout(p=0.5, inplace=False)
  )
  (decoder): DecoderWithAttention(
    (attention): Attention(
      (attn): Linear(in_features=1024, out_features=512, bias=True)
      (v): Linear(in_features=512, out_features=1, bias=True)
    )
    (embedding): Embedding(10104, 256)
    (rnn): GRU(768, 512, num_layers=2, dropout=0.5)
    (out): Linear(in_features=1280, out_features=10104, bias=True)
    (dropout): Dropout(p=0.5, inplace=False)
  )
)

```

```

1 def count_parameters(model):
2     return sum(p.numel() for p in model.parameters() if p.requires_grad)
3
4 print(f'The model has {count_parameters(model):,} trainable parameters')

```

The model has 25,846,905 trainable parameters

Let's take a look at our network quality:

▼ Bleu

[link](#)

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}.$$

Then,

$$\text{BLEU} = \text{BP} \cdot \exp \left(\sum_{n=1}^N w_n \log p_n \right).$$

The ranking behavior is more immediately apparent in the log domain,

$$\log \text{BLEU} = \min\left(1 - \frac{r}{c}, 0\right) + \sum_{n=1}^N w_n \log p_n.$$

In our baseline, we use $N = 4$ and uniform weights

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

1 def cut_on_eos(tokens_iter):
2     for token in tokens_iter:
3         if token == '<eos>':
4             break
5         yield token
6
7 def remove_tech_tokens(tokens_iter, tokens_to_remove=['<sos>', '<unk>', '<pad>']):
8     return [x for x in tokens_iter if x not in tokens_to_remove]
9
10 def generate_translation(src, trg, model, TRG_vocab, SRC_vocab):
11     model.eval()
12
13     output = model(src, trg, 0) #turn off teacher forcing
14     output = output[1:].argmax(-1)
15
16     source = remove_tech_tokens(cut_on_eos([SRC_vocab.itos[x] for x in list(src[:,0].cpu().numpy())]))

```

```

17 original = remove_tech_tokens(cut_on_eos([TRG_vocab.itos[x] for x in list(trg[:,0].
18 generated = remove_tech_tokens(cut_on_eos([TRG_vocab.itos[x] for x in list(output[:
19
20 print('Source: {}'.format(' '.join(source[:::-1])))
21 print('Original: {}'.format(' '.join(original)))
22 print('Generated: {}'.format(' '.join(generated)))
23 print()
24
25 def get_text(x, vocab):
26     generated = remove_tech_tokens(cut_on_eos([vocab.itos[elem] for elem in list(x)]))
27     return generated

```

```

1 from nltk.translate.bleu_score import corpus_bleu
2
3 # """ Estimates corpora-level BLEU score of model's translations given inp and refe
4 # translations, _ = model.translate_lines(inp_lines, **flags)
5 # # Note: if you experience out-of-memory error, split input lines into batches and
6 # return corpus_bleu([[ref] for ref in out_lines], translations) * 100
7
8 def train(model, iterator, optimizer, criterion, clip, train_history=None, valid_histor
9     model.train()
10
11 epoch_loss = 0
12 history = []
13 for i, batch in enumerate(tqdm.notebook.tqdm(iterator)):
14
15     src = batch.src
16     trg = batch.trg
17
18     optimizer.zero_grad()
19
20     output = model(src, trg)
21
22     #trg = [trg sent len, batch size]

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

26     trg = trg[1:].view(-1)
27
28     #trg = [(trg sent len - 1) * batch size]
29     #output = [(trg sent len - 1) * batch size, output dim]
30
31     loss = criterion(output, trg)
32
33     loss.backward()
34
35     # Let's clip the gradient
36     torch.nn.utils.clip_grad_norm_(model.parameters(), clip)
37
38     optimizer.step()
39
40     epoch_loss += loss.item()
41
42     history.append(loss.cpu().data.numpy())

```

```

43
44     return history, (epoch_loss / len(iterator))
45
46 def evaluate(model, iterator, criterion):
47
48     model.eval()
49
50     epoch_loss = 0
51
52     history = []
53
54     with torch.no_grad():
55         for i, batch in enumerate(tqdm.notebook.tqdm(iterator)):
56
57             src = batch.src
58             trg = batch.trg
59
60             output = model(src, trg, 0) #turn off teacher forcing
61
62             #trg = [trg sent len, batch size]
63             #output = [trg sent len, batch size, output dim]
64
65             output = output[1:].view(-1, OUTPUT_DIM)
66             trg = trg[1:].view(-1)
67
68             #trg = [(trg sent len - 1) * batch size]
69             #output = [(trg sent len - 1) * batch size, output dim]
70
71             loss = criterion(output, trg)
72
73             epoch_loss += loss.item()
74
75     return epoch_loss / len(iterator)
76
77

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

80     generated_text = []
81
82     model.eval()
83     with torch.no_grad():
84         for i, batch in enumerate(tqdm.notebook.tqdm(iterator)):
85
86             src = batch.src
87             trg = batch.trg
88
89             output = model(src, trg, 0) #turn off teacher forcing
90
91             #trg = [trg sent len, batch size]
92             #output = [trg sent len, batch size, output dim]
93
94             output = output[1:].argmax(-1)
95
96             original_text.extend([get_text(x, TRG_vocab) for x in trg.cpu().numpy().T])
97             generated_text.extend([get_text(x, TRG_vocab) for x in output.detach().cpu(

```

```

98         # original_text = flatten(original_text)
99         # generated_text = flatten(generated_text)
100     return corpus_bleu([[text] for text in original_text], generated_text) * 100
101
102
103
104 def epoch_time(start_time, end_time):
105     elapsed_time = end_time - start_time
106     elapsed_mins = int(elapsed_time / 60)
107     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
108     return elapsed_mins, elapsed_secs

```

```

1 PAD_IDX = TRG.vocab.stoi['<pad>']
2 optimizer = optim.Adam(model.parameters()) # ORIGINAL: SGD together with Adadelta
3 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(optimizer, 'min')
4 criterion = nn.CrossEntropyLoss(ignore_index = PAD_IDX)

```

```

1 import matplotlib
2 matplotlib.rcParams.update({'figure.figsize': (16, 12), 'font.size': 14})
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 from IPython.display import clear_output

```

```

1 FROM_PRETRAINED = False

```

```

1 def init_weights(m):
2     for name, param in m.named_parameters():
3         nn.init.uniform_(param, -0.08, 0.08)
4
5 #model.apply(init_weights)
6
7 if FROM_PRETRAINED:

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

11     model.load_state_dict(torch.load('/content/drive/MyDrive/Colab Notebooks/Attent
12     print('Loaded pre-trained model.')
13 except FileNotFoundError as e:
14     print(e)
15     model.apply(init_weights)
16     print('\nInitialized the weights randomly: Uniform distrubution (-0.08, 0.08).')
17
18 else:
19     model.apply(init_weights)
20     print('\nInitialized the weights randomly: Uniform distrubution (-0.08, 0.08).')

```

Initialized the weights randomly: Uniform distrubution (-0.08, 0.08).

```

1 TRAIN = True

```

```

1 import tqdm
2
3 if TRAIN:
4     train_history = []
5     valid_history = []
6     valid_bleu_history = []
7
8     N_EPOCHS = 30
9     PATIENCE = 5
10    CLIP = 5
11
12    # better to optimise for BLEU on validation set if that's what we care about on
13    # the test set
14    # best_valid_loss = float('inf')
15
16    # allows to keep the best model saved and only substituted if better valid bleu is
17    print('Calculating valid bleu of the best checkpoint.')
18    best_valid_bleu = evaluate_bleu(model, valid_iterator, TRG.vocab)
19    print(f'Best valid bleu achieved: {best_valid_bleu:.3}')
20    print('\n', '-'*80, '\n')
21
22    best_epoch = 0
23    current_patience = 0
24
25    for epoch in range(N_EPOCHS):
26        print(f'Epoch: {epoch+1:02}')
27
28        start_time = time.time()
29
30        print('Calculating train_loss')
31        epochs_history, train_loss = train(model, train_iterator, optimizer, criterion,
32        print('Calculating valid_loss')
33        valid_loss = evaluate(model, valid_iterator, criterion)
34        print('Calculating valid_bleu')
35        valid_bleu = evaluate_bleu(model, valid_iterator, TRG.vocab)
36
37        scheduler.step(valid_loss)
38
39        end_time = time.time()
40
41        epoch_mins, epoch_secs = epoch_time(start_time, end_time)
42
43        # REMEMBER TO CHECK THE SIGN
44        if valid_bleu > best_valid_bleu:
45            # record
46            best_valid_bleu = valid_bleu
47            current_patience = 0
48            best_epoch = epoch
49            # save
50            torch.save(model.state_dict(), '/content/drive/MyDrive/Colab Notebooks/Atte
51        else:
52            current_patience += 1
53
54        train_history.append(train_loss)

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

```

56     valid_history.append(valid_loss)
57     valid_bleu_history.append(valid_bleu)
58
59     # plot once every epoch
60     fig, ax = plt.subplots(nrows=1, ncols=3, figsize=(13.5, 6))
61     ax[0].plot(epochs_history, label='train loss')
62     ax[0].set_xlabel('Batch')
63     ax[0].set_title('Train loss')
64     if train_history is not None:
65         ax[1].plot(train_history, label='train loss')
66         ax[1].set_xlabel('Epoch')
67     if valid_history is not None:
68         ax[1].plot(valid_history, label='valid loss')
69     if valid_bleu_history is not None:
70         ax[2].plot(valid_bleu_history, label='valid BLEU')
71         ax[2].set_xlabel('Epoch')
72     plt.legend()
73     plt.show()
74     # print once every epoch
75     print(f'Epoch: {epoch+1:02} | Time: {epoch_mins}m {epoch_secs}s')
76     print(f'\tTrain Loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):7.3f}')
77     print(f'\tVal. Loss: {valid_loss:.3f} | Val. PPL: {math.exp(valid_loss):7.3f}')
78     print(f'\tVal. BLEU: {valid_bleu:.3f}')
79
80     # break if reached the patience
81     if current_patience > PATIENCE:
82         print(f"No improvement for {PATIENCE} epochs.")
83         break
84
85     print('\n', '-'*80, '\n')
86
87     print(f'Best valid bleu = {best_valid_bleu:.3f} achieved at epoch {best_epoch+1:02}')
88 else:
89     print('Using a pre-trained model w/o further training.')

```

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu ✕

Calculating valid bleu of the best checkpoint.

100% 20/20 [00:05<00:00, 2.10it/s]

Best valid bleu achieved: 8.0

Epoch: 01

Calculating train_loss

/usr/local/lib/python3.7/dist-packages/nltk/translate/bleu_score.py:490: UserWarning
Corpus/Sentence contains 0 counts of 2-gram overlaps.

BLEU scores might be undesirable; use SmoothingFunction().

warnings.warn(_msg)

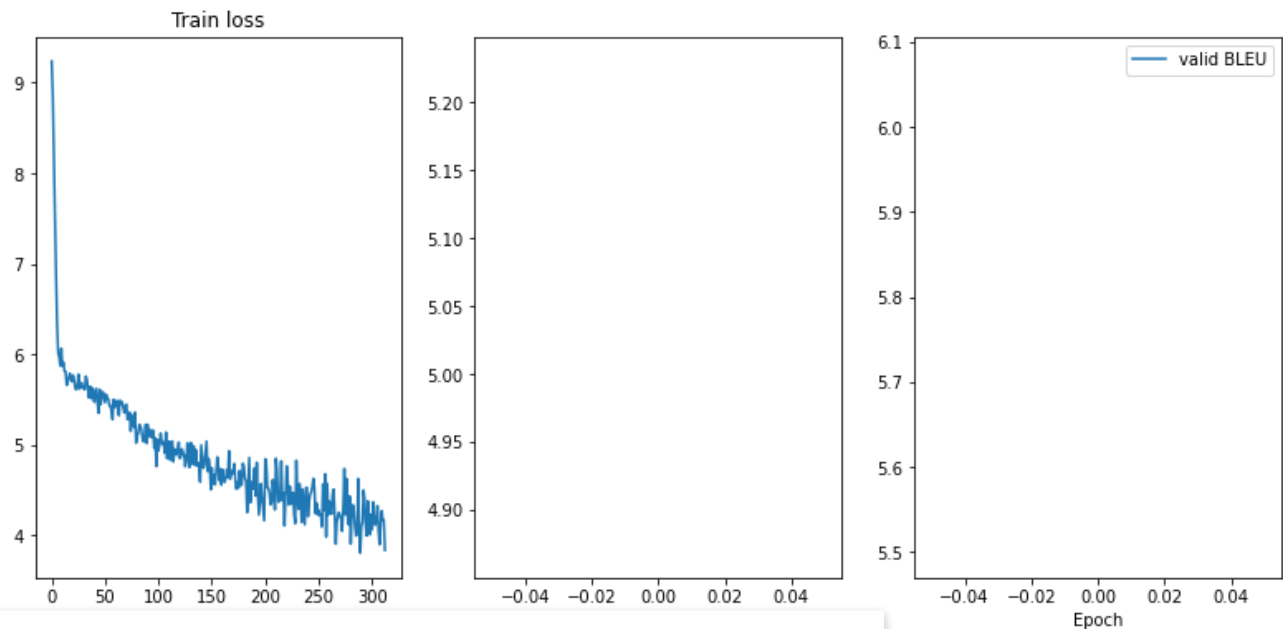
100% 313/313 [08:11<00:00, 1.41s/it]

Calculating valid_loss

100% 20/20 [00:05<00:00, 2.10it/s]

Calculating valid_bleu

100% 20/20 [00:05<00:00, 2.11it/s]



To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu

Val. Loss: 5.230 | Val. PPL: 186.800
Val. BLEU: 5.788

Epoch: 02

Calculating train_loss

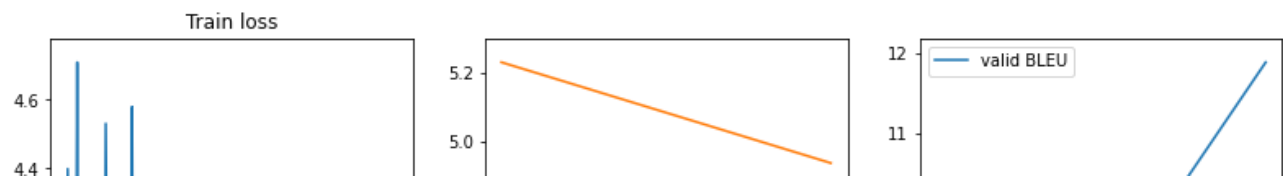
100% 313/313 [08:11<00:00, 1.61s/it]

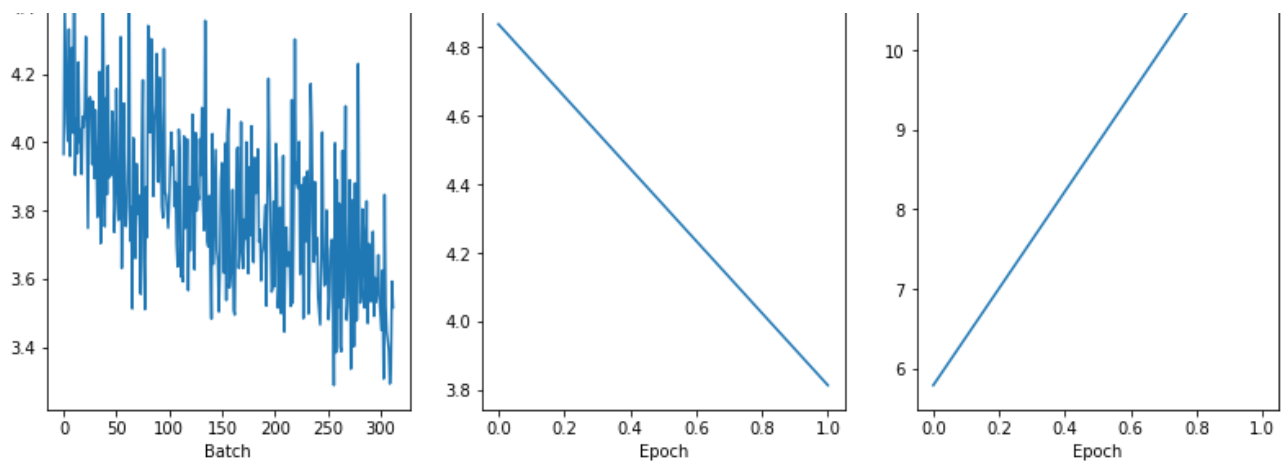
Calculating valid_loss

100% 20/20 [00:05<00:00, 2.13it/s]

Calculating valid_bleu

100% 20/20 [00:05<00:00, 2.12it/s]





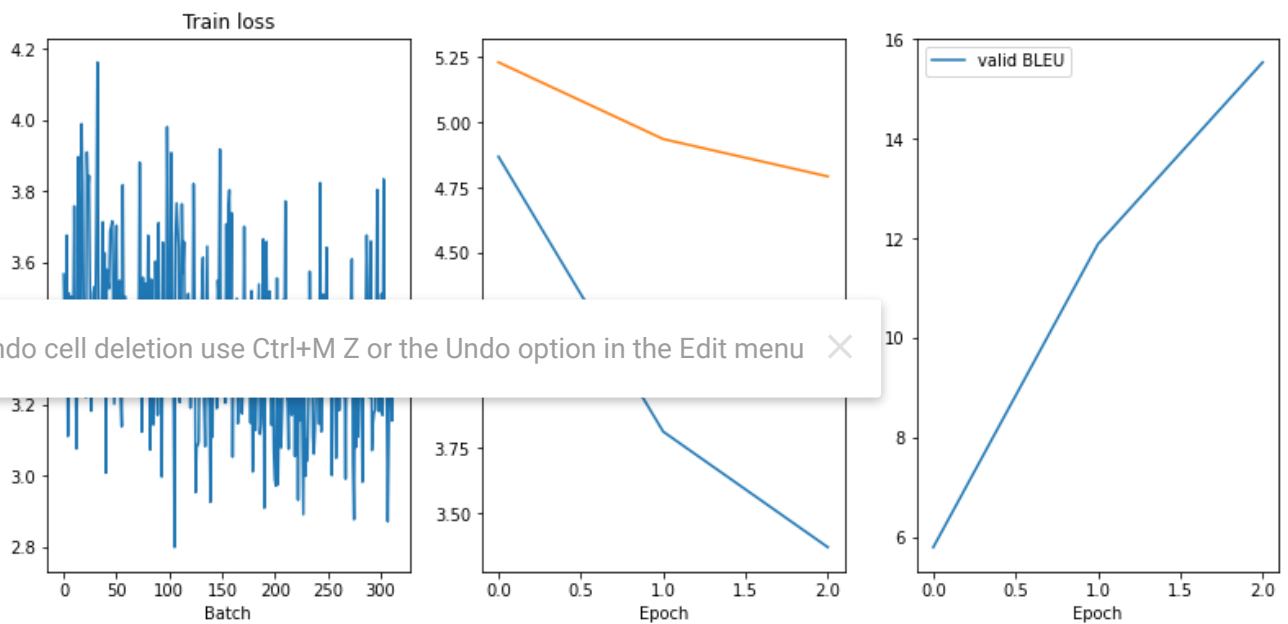
Epoch: 02 | Time: 8m 22s
 Train Loss: 3.813 | Train PPL: 45.308
 Val. Loss: 4.936 | Val. PPL: 139.160
 Val. BLEU: 11.885

Epoch: 03
 Calculating train_loss
 100%
 Calculating valid_loss
 100%
 Calculating valid_bleu
 100%

313/313 [08:13<00:00, 1.51s/it]

20/20 [00:05<00:00, 2.13it/s]

20/20 [00:05<00:00, 2.11it/s]



To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu

Epoch: 03 | Time: 8m 24s
 Train Loss: 3.370 | Train PPL: 29.072
 Val. Loss: 4.792 | Val. PPL: 120.568
 Val. BLEU: 15.538

Epoch: 04
 Calculating train_loss
 100%
 Calculating valid_loss
 100%

313/313 [08:13<00:00, 1.42s/it]

20/20 [00:05<00:00, 2.08it/s]

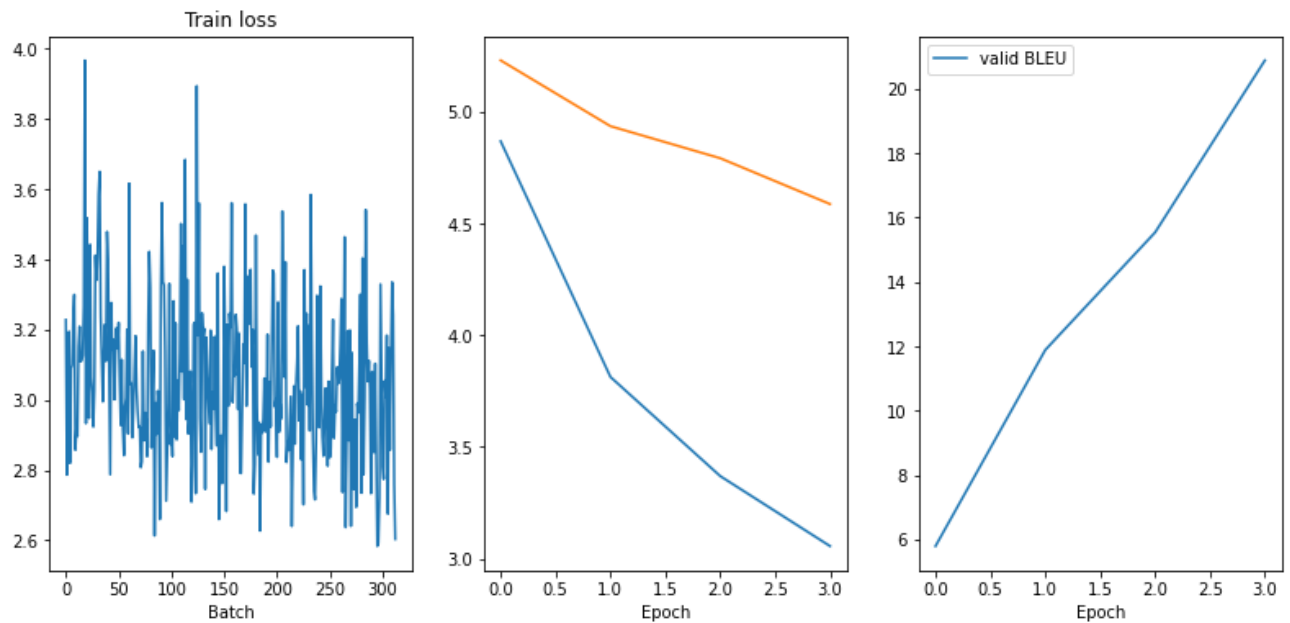
100%

20/20 [00:05<00:00, 2.00it/s]

Calculating valid_bleu

100%

20/20 [00:05<00:00, 2.11it/s]



Epoch: 04 | Time: 8m 24s

Train Loss: 3.057 | Train PPL: 21.254

Val. Loss: 4.587 | Val. PPL: 98.204

Val. BLEU: 20.886

Epoch: 05

Calculating train_loss

100%

313/313 [08:15<00:00, 1.63s/it]

Calculating valid_loss

100%

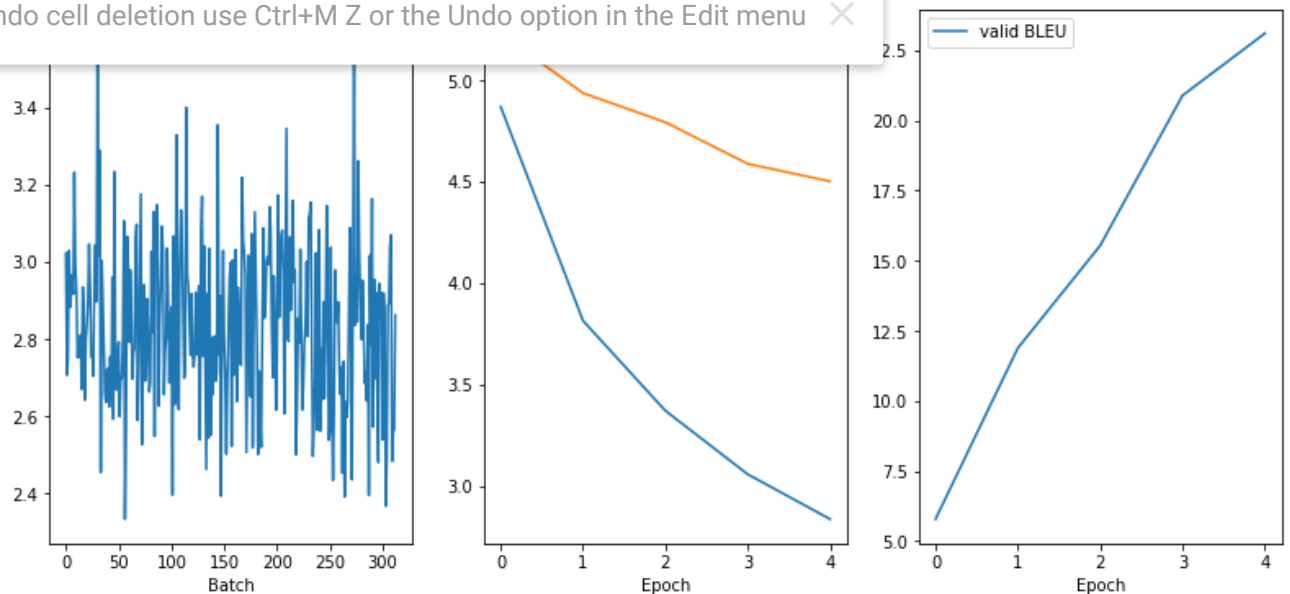
20/20 [00:05<00:00, 2.10it/s]

Calculating valid_bleu

100%

20/20 [00:05<00:00, 2.10it/s]

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu



Epoch: 05 | Time: 8m 26s

Train Loss: 2.835 | Train PPL: 17.037

Val. Loss: 4.500 | Val. PPL: 90.056

Val. BLEU: 23.100

Epoch: 06

Calculating train_loss

100%

313/313 [08:16<00:00, 1.26s/it]

Calculating valid_loss

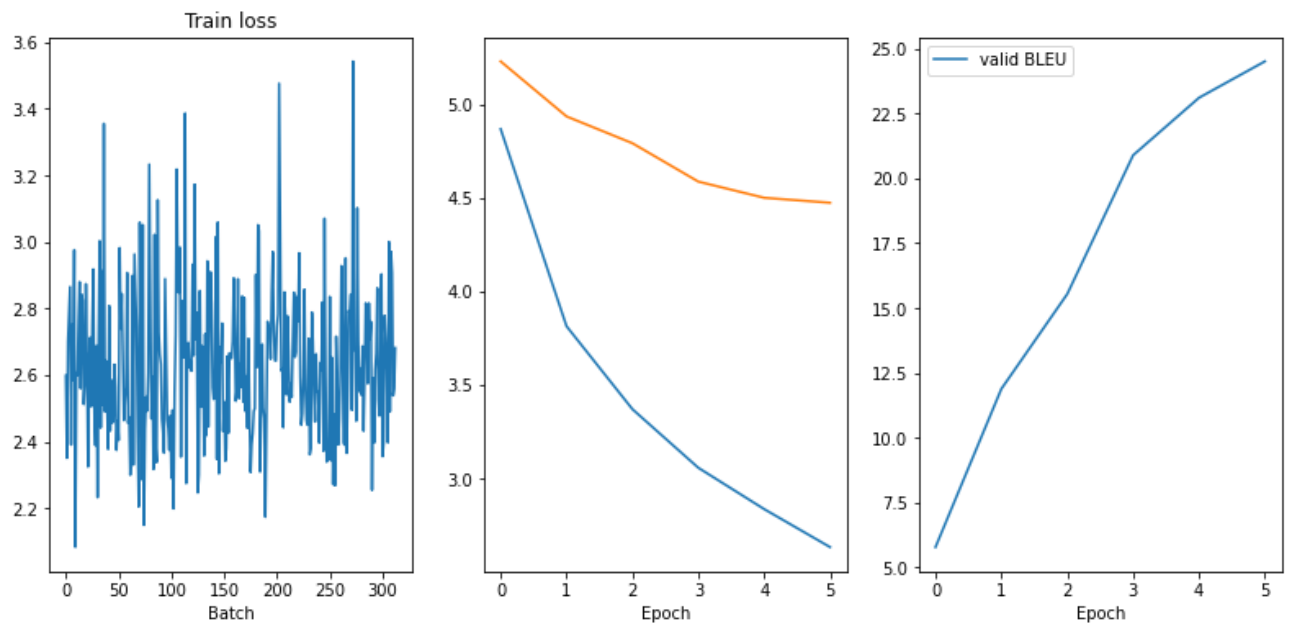
100%

20/20 [00:05<00:00, 2.11it/s]

Calculating valid_bleu

100%

20/20 [00:05<00:00, 2.11it/s]



Epoch: 06 | Time: 8m 27s

Train Loss: 2.632 | Train PPL: 13.903

Val. Loss: 4.474 | Val. PPL: 87.702

Val. BLEU: 24.502

Epoch: 07

To undo cell deletion use Ctrl+M Z or the Undo option in the Edit menu

20/20 [00:05<00:00, 1.68s/it]

Calculating valid_loss

100%

20/20 [00:05<00:00, 2.09it/s]

Calculating valid_bleu

100%

20/20 [00:05<00:00, 2.08it/s]

