

This can be run [run on Google Colab using this link \(https://colab.research.google.com/github/CS7150/CS7150-Homework\\_4/blob/main/Assignment\\_4\\_Transformers.ipynb\)](https://colab.research.google.com/github/CS7150/CS7150-Homework_4/blob/main/Assignment_4_Transformers.ipynb)

## Dependencies

```
In [1]: 1 |pip install -U spacy==3.6.0
2 |python -m spacy download en_core_web_sm
3 |python -m spacy download de_core_news_sm
4 |pip install torchdata
5 |pip install -U torchtext
6 |pip install portalocker>=2.0.0
7 |pip install seaborn
8 |python -m tensorflow/compiler/xla/stream_executor/cuda/cuda_llvmjit.py --enable-to-register-cudnn-factory: Attempting to register factory for
or plugin cuDNN when one has already been registered
2023-11-22 01:37:24.078054: E tensorflow/compiler/xla/stream_executor/cuda/cuda_fft.cc:609] Unable to register cuFFT factory: Attempting to register factory fo
r plugin cuFFT when one has already been registered
2023-11-22 01:37:24.078089: E tensorflow/compiler/xla/stream_executor/cuda/cuda_blas.cc:1518] Unable to register cuBLAS factory: Attempting to register factory
for plugin cuBLAS when one has already been registered
2023-11-22 01:37:24.086305: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in perfo
rmance-critical operations.
To enable the following instructions: AVX2 AVX512F FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2023-11-22 01:37:25.256336: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
2023-11-22 01:37:26.787845: I tensorflow/compiler/xla/stream_executor/cuda/cuda_gpu_executor.cc:894] successful NUMA node read from SysFS had negative value (-
1), but there must be at least one NUMA node, so returning NUMA node zero. See more at https://github.com/torvalds/linux/blob/v6.0/Documentation/ABI/testing/sysfs-bus-pci#L344-L355
2023-11-22 01:37:26.788272: I tensorflow/compiler/xla/stream_executor/cuda/cuda_gpu_executor.cc:894] successful NUMA node read from SysFS had negative value (-
1), but there must be at least one NUMA node, so returning NUMA node zero. See more at https://github.com/torvalds/linux/blob/v6.0/Documentation/ABI/testing/sysfs-bus-pci#L344-L355
2023-11-22 01:37:26.788448: I tensorflow/compiler/xla/stream_executor/cuda/cuda_gpu_executor.cc:894] successful NUMA node read from SysFS had negative value (-
1), but there must be at least one NUMA node, so returning NUMA node zero. See more at https://github.com/torvalds/linux/blob/v6.0/Documentation/ABI/testing/sysfs-bus-pci#L344-L355
Collecting de-core-news-sm==3.6.0
```

## Transformer Assignment

## Overview

In this assignment, you will be trying your hand at understanding transformers, their architecture, and their difference in-terms of basic RNNs. The assignment is divided in 2 sections.

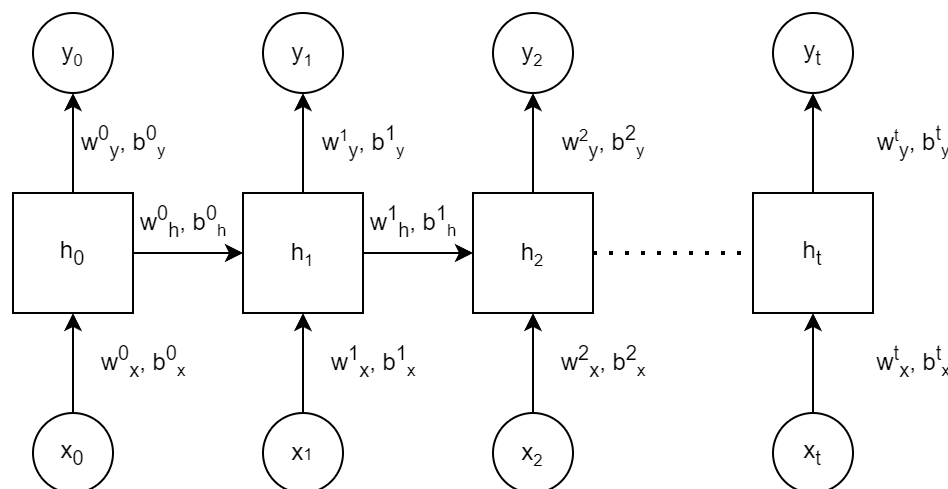
- Section 1:  
You will be implementing a basic RNN cell, RNN Class and an RNN Classifier
- Section 2:  
You will be implementing a Transformer based Text classifier using components such as Multi-head Attention Module, Positional Encoding Module and Encoder
- Section 3:  
In order to experiment with Decoders for a Transformer, we will be implementing a Transformer Based Machine Translation class using modules of Section 2, a Decoder, Attention Masks and Seq-Seq Module

```
In [2]: 1 import math
2 import torch
3 import time
4
5 import numpy as np
6 import torch.nn as nn
7 import torch.nn.functional as F
8 import pandas as pd
9 import seaborn as sns
10
11 from torch.utils.data import DataLoader
12 from torchtext.datasets import AG_NEWS
13 from torch.utils.data.dataset import random_split
14 from torchtext.data.functional import map_style_dataset
15 from torchtext.data.utils import get_tokenizer
16 from torchtext.vocab import build_vocab_from_iterator
17 from torchtext.datasets import Multi30k
18 from typing import Iterable, List
```

## Section 1: Recurrent Neural Networks (RNN)

Each RNN Cell should contain 2 components: an Input Unit and a Hidden Unit. The Hidden state is the part of the RNN that remembers context about previous data present in the sequence. The current time step's hidden state is calculated using information of the previous time steps hidden state and the current input. This process helps to retain information on what the model saw in the previous time step when processing the current time steps information.

RNNs will look and function as follows.



The hidden state any given time  $t$  is given by,

$$\begin{aligned} input_t &= (x_t \cdot W_x^t + b_x^t) \\ prev\_state &= (h_{t-1} \cdot W_h^t + b_h^t) \\ h_t &= \tanh(input_t + prev\_state) \end{aligned}$$

The output at any give time  $t$  is given by,

$$y_t = h_t \cdot W_y^t + b_y$$

Note: All the connections in RNN have weights and biases.

Your job is to implement the formulae above.

## 1.1 A Single RNN Cell

```
In [3]: 1 class RNNCell(torch.nn.Module):
2         """
3         RNNCell is a single cell that takes x_t and h_{t-1} as input and outputs h_t.
4         """
5         def __init__(self, input_dim: int, hidden_dim: int):
6             """
7             Constructor of RNNCell.
8
9             Inputs:
10            - input_dim: Dimension of the input x_t
11            - hidden_dim: Dimension of the hidden state h_{t-1} and h_t
12            """
13
14            # We always need to do this step to properly implement the constructor
15            super(RNNCell, self).__init__()
16
17            #self.linear_x, self.linear_h, self.non_linear = None, None, None
18
19            #####
20            # TODO:
21            # 1. Define the linear transformation layers for the attributes
22            #    (set to None above) to correspond to the W_x and W_h in the formulae.
23            #    Remember to include bias in the linear layers.
24            #    (Refer to nn.Linear documentation https://pytorch.org/docs/stable/generated/torch.nn.Linear.html)
25            # 2. Define the non_linear layer. (You can use tanh as describe bove).
26            #####
27            self.linear_x = torch.nn.Linear(input_dim, hidden_dim) #x
28            self.linear_h = torch.nn.Linear(hidden_dim, hidden_dim) #h
29            self.non_linear = torch.nn.Tanh() #nonlinear
30
31            #####
32            #                      END OF YOUR CODE                      #
33            #####
34
35            def forward(self, x_cur: torch.Tensor, h_prev: torch.Tensor):
36                """
37                Compute h_t given x_t and h_{t-1}.
38
39                Inputs:
40                - x_cur: x_t, a tensor with the same of BxC, where B is the batch size and
41                  C is the channel dimension.
42                - h_prev: h_{t-1}, a tensor with the same of BxH, where H is the channel
43                  dimension.
44                """
45                #h_cur = None
46                #####
47                # TODO: Run the linear transformation layers to compute x_t and consume
48                #    h_{t-1}
49                # go non-linear layer.
50                #####
51                x_transformed = self.linear_x(x_cur) #x_transformed
52                h_transformed = self.linear_h(h_prev) #h_transformed
53                h_cur = self.non_linear(x_transformed + h_transformed) #current_h
54                #####
55                #                      END OF YOUR CODE                      #
56                #####
57                return h_cur
```

```
In [4]: 1 # Let's run a sanity check of your model
2 x = torch.randn((2, 8)) # Input Dim
3 h = torch.randn((2, 16)) # Hidden Dim
4
5 model = RNNCell(8, 16)
6 y = model(x, h)
7 assert len(y.shape) == 2 and y.shape[0] == 2 and y.shape[1] == 16
8 print(y.shape)
```

torch.Size([2, 16])



## 1.2 RNN Layer

```
In [5]: 1 class RNN(torch.nn.Module):
2         """
3         RNN is a single-layer (stack) RNN by connecting multiple RNNCell together in a single
4         direction, where the input sequence is processed from left to right.
5         """
6         def __init__(self, input_dim: int, hidden_dim: int):
7             """
8             Constructor of the RNN module.
9
10            Inputs:
11            - input_dim: Dimension of the input x_t
12            - hidden_dim: Dimension of the hidden state h_{t-1} and h_t
13            """
14            super(RNN, self).__init__()
15            self.hidden_dim = hidden_dim
16
17            #####
18            # TODO: Define the RNNCell. #
19            #####
20
21            self.rnn_cell = RNNCell(input_dim, hidden_dim) #rnn_cell
22            #####
23            #                               END OF YOUR CODE                               #
24            #####
25
26        def forward(self, x: torch.Tensor):
27            """
28            Compute the hidden representations for every token in the input sequence.
29
30            Input:
31            - x: A tensor with the shape of BxLxC, where B is the batch size, L is the sequence
32              length, and C is the channel dimension
33
34            Return:
35            - h: A tensor with the shape of BxLxH, where H is the hidden dimension of RNNCell
36            """
37            b = x.shape[0]
38            seq_len = x.shape[1]
39
40            # initialize the hidden dimension
41            init_h = x.new_zeros((b, self.hidden_dim))
42            #h = None
43            h = []
44            #####
45            # TODO: Compute the hidden representation for every token in the input #
46            # from left to right as per the formula stated above #
47            #####
48            # Computing the hidden representation for every token
49            prev_hidden = init_h
50            for t in range(seq_len):
51                cur_hidden = self.rnn_cell(x[:, t, :], prev_hidden)
52                h.append(cur_hidden.unsqueeze(1))
53                prev_hidden = cur_hidden
54
55            h = torch.cat(h, dim=1)
56
57            #####
58            #                               END OF YOUR CODE                               #
59            #####
60
61            return h
```

```
In [6]: 1 # Let's run a sanity check of your model
2 x = torch.randn((2, 10, 8))
3 model = RNN(8, 16)
4 y = model(x)
5 assert len(y.shape) == 3
6 for dim, dim_gt in zip(y.shape, [2, 10, 16]):
7     assert dim == dim_gt
8 print(y.shape)
```

torch.Size([2, 10, 16])



## 1.3 RNN Classifier

```

In [7]: 1 h_tracker = {}
        2
        3 class RNNClassifier(nn.Module):
        4     """
        5     A RNN-based classifier for text classification. It first converts tokens into word embeddings.
        6     And then feeds the embeddings into a RNN, where the hidden representations of all tokens are
        7     then averaged to get a single embedding of the sentence. It will be used as input to a linear
        8     classifier.
        9     """
        10    def __init__(self,
        11                vocab_size: int, embed_dim: int, rnn_hidden_dim: int, num_class: int, pad_token: int
        12                ):
        13        """
        14        Constructor.
        15
        16        Inputs:
        17        - vocab_size: Vocabulary size, indicating how many tokens we have in total.
        18        - embed_dim: The dimension of word embeddings
        19        - rnn_hidden_dim: The hidden dimension of the RNN.
        20        - num_class: Number of classes.
        21        - pad_token: The index of the padding token.
        22        """
        23        super(RNNClassifier, self).__init__()
        24
        25        # word embedding layer
        26        self.embedding = nn.Embedding(vocab_size, embed_dim, padding_idx=pad_token) #embedding
        27
        28        #####
        29        # TODO: Define the RNN and the classification layer. #
        30        #####
        31
        32        self.rnn = nn.RNN(embed_dim, rnn_hidden_dim) #rnn
        33
        34        self.fc = nn.Linear(rnn_hidden_dim, num_class) #fc
        35
        36        self.init_weights()
        37
        38        #####
        39        # END OF YOUR CODE #
        40        #####
        41
        42    def init_weights(self):
        43        initalize = 0.5
        44        self.embedding.weight.data.uniform_(-initalize, initalize)
        45        self.fc.weight.data.uniform_(-initalize, initalize)
        46        self.fc.bias.data.zero_()
        47
        48    def forward(self, text):
        49        """
        50        Get classification scores (logits) of the input.
        51
        52        Input:
        53        - text: Tensor with the shape of BxLxC.
        54
        55        Return:
        56        - logits: Tensor with the shape of BxK, where K is the number of classes
        57        """
        58
        59        # get word embeddings
        60        embedded = self.embedding(text)
        61
        62        #####
        63        # TODO: Compute logits of the input. #
        64        #####
        65
        66        rnn_output, _ = self.rnn(embedded) # RNN layer
        67
        68        avg_pool = torch.mean(rnn_output, dim=1) # Average pooling over time steps
        69
        70        logits = self.fc(avg_pool) # Classification layer
        71        #####
        72        # END OF YOUR CODE #
        73        #####
        74
        75        return logits

```

```

In [8]: 1 # Sanity check!!!
        2 vocab_size = 10
        3 embed_dim = 16
        4 rnn_hidden_dim = 32
        5 num_class = 4
        6
        7 x = torch.arange(vocab_size).view(1, -1)
        8 x = torch.cat((x, x), dim=0)
        9 print('x.shape: {}'.format(x.shape))
        10 model = RNNClassifier(vocab_size, embed_dim, rnn_hidden_dim, num_class, 0)
        11 y = model(x)
        12 assert len(y.shape) == 2 and y.shape[0] == 2 and y.shape[1] == num_class
        13 print(y.shape)

```

x.shape: torch.Size([2, 10])  
torch.Size([2, 4])

## Data Loader

```

In [9]: 1 # check here for details https://github.com/pytorch/text/blob/main/torchtext/data/utils.py#L52-L166
2 from torchtext.data.utils import get_tokenizer
3 # check here for details https://github.com/pytorch/text/blob/main/torchtext/vocab/vocab_factory.py#L65-L113
4 from torchtext.vocab import build_vocab_from_iterator
5 # Documentation of DataLoader https://pytorch.org/docs/stable/data.html#torch.utils.data.DataLoader
6 from torch.utils.data import DataLoader
7
8 # A tokenizer splits a input setence into a set of tokens, including those puncuation
9 # For example
10 # >>> tokens = tokenizer("You can now install TorchText using pip!")
11 # >>> tokens
12 # >>> ['you', 'can', 'now', 'install', 'torchtext', 'using', 'pip', '!']
13 tokenizer = get_tokenizer('basic_english')
14
15 train_iter = AG_NEWS(split='train')
16
17 def yield_tokens(data_iter):
18     for _, text in data_iter:
19         yield tokenizer(text)
20
21 # Creates a vocab object which maps tokens to indices
22 # Check here for details https://github.com/pytorch/text/blob/main/torchtext/vocab/vocab.py
23 vocab = build_vocab_from_iterator(yield_tokens(train_iter), specials=["<unk>"])
24
25 # The specified token will be returned when a out-of-vocabulary token is queried.
26 vocab.set_default_index(vocab["<unk>"])
27
28 text_pipeline = lambda x: vocab(tokenizer(x))
29 label_pipeline = lambda x: int(x) - 1
30
31 # The padding token we need to use
32 # The returned indices are always in an array
33 PAD_TOKEN = vocab(tokenizer('<pad>'))
34 assert len(PAD_TOKEN) == 1
35 PAD_TOKEN = PAD_TOKEN[0]
36
37 # Merges a list of samples to form a mini-batch of Tensor(s)
38 def collate_batch(batch):
39     """
40     Input:
41     - batch: A list of data in a mini batch, where the length denotes the batch size.
42       The actual context depends on a particular dataset. In our case, each position
43       contains a label and a Tensor (tokens in a sentence).
44
45     Returns:
46     - batched_label: A Tensor with the shape of (B,)
47     - batched_text: A Tensor with the shape of (B, L, C), where L is the sequence length
48       and C is the channel dimension
49     """
50     label_list, text_list, text_len_list = [], [], []
51     for (_label, _text) in batch:
52         label_list.append(label_pipeline(_label))
53         processed_text = torch.tensor(text_pipeline(_text), dtype=torch.int64)
54         text_list.append(processed_text)
55         text_len_list.append(processed_text.size(0))
56
57     #####
58     # TODO: Pad the text tensor in the mini batch so that they have the same #
59     # length. Specifically, you need to calculate the maximum length in the #
60     # batch and then add the token PAD_TOKEN to the end of those #
61     # shorter sentences. (Try printing a few data points to understand why) #
62     #####
63     max_len = max(len(text) for text in text_list) #getting max length
64     padded_text = [torch.cat((text, torch.tensor([PAD_TOKEN] * (max_len - len(text))), dtype=torch.int64)) for text in text_list] #padded text
65     batched_text = torch.stack(padded_text, dim=0) # batched_text
66     batched_label = torch.tensor(label_list, dtype=torch.int64) # batched_label
67     #####
68     # END OF YOUR CODE #
69     #####
70
71     return batched_label.long(), batched_text.long()
72
73 # Now, let's check what the batched data looks like
74 train_iter = AG_NEWS(split='train')
75 dataloader = DataLoader(train_iter, batch_size=8, shuffle=False, collate_fn=collate_batch)
76 for idx, (label, data) in enumerate(dataloader):
77     if idx > 0:
78         break
79     print('label.shape: {}'.format(label.shape))
80     print('label: {}'.format(label))
81     print('data.shape: {}'.format(data.shape))
82
83 label.shape: torch.Size([8])
84 label: tensor([2, 2, 2, 2, 2, 2, 2, 2])
85 data.shape: torch.Size([8, 49])

```

```

In [10]: 1 labels=set()
2 labels.update([entry[0] for entry in AG_NEWS(root="data")[0]])
3 print(labels)

{1, 2, 3, 4}

```

## 1.4 Train & Evaluate Module

```
In [11]: 1 # logits_tracker = {}
2 def train(model, dataloader, loss_func, device, grad_norm_clip, optimizer):
3     model.train()
4     total_acc, total_count = 0, 0
5     log_interval = 500
6     start_time = time.time()
7     global logits_tracker
8
9     for idx, (label, text) in enumerate(dataloader):
10        label = label.to(device)
11        text = text.to(device)
12        optimizer.zero_grad()
13
14        #####
15        # TODO: compute the logits of the input, get the loss, and do the #
16        # gradient backpropagation.
17        #####
18
19        logits = model(text) #model logits
20        loss = loss_func(logits, label) # Calculating the loss
21        loss.backward() # Backpropagating the gradients
22
23        #####
24        #                               END OF YOUR CODE                               #
25        #####
26
27        torch.nn.utils.clip_grad_norm_(model.parameters(), grad_norm_clip)
28        optimizer.step()
29        total_acc += (logits.argmax(1) == label).sum().item()
30        total_count += label.size(0)
31        if idx % log_interval == 0 and idx > 0:
32            elapsed = time.time() - start_time
33            print('| epoch {:3d} | {:5d}/{:5d} batches |
34                  | accuracy {:.3f}'.format(epoch, idx, len(dataloader),
35                                            total_acc/total_count))
36            total_acc, total_count = 0, 0
37            start_time = time.time()
38
39 def evaluate(model, dataloader, loss_func, device):
40     model.eval()
41     total_acc, total_count = 0, 0
42
43     with torch.no_grad():
44         for idx, (label, text) in enumerate(dataloader):
45             label = label.to(device)
46             text = text.to(device)
47
48             #####
49             # TODO: compute the logits of the input, get the loss. #
50             #####
51             logits = model(text) # Computing logits of the input
52             loss = loss_func(logits, label) # Calculating the loss
53             #####
54             #                               END OF YOUR CODE                               #
55             #####
56
57             total_acc += (logits.argmax(1) == label).sum().item()
58             total_count += label.size(0)
59     return total_acc/total_count
```

```
In [12]: 1
2 assert torch.cuda.is_available(), "Please connect to the GPU instance if working on Colab or configure the environment for Torch using GPU (Comment this line
3 # device = 'cuda'
4 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
5
6 # Hyper parameters
7 epochs = 3 # epoch
8 lr = 0.0005 # learning rate
9 batch_size = 64 # batch size for training
10 word_embed_dim = 64
11 rnn_hidden_dim = 96
12
13 train_iter = AG_NEWS(split='train')
14 num_class = len(set([label for (label, text) in train_iter]))
15 vocab_size = len(vocab)
16
17 #####
18 # TODO: Define the classifier and loss function.
19 #####
20
21 model = RNNClassifier(vocab_size, word_embed_dim, rnn_hidden_dim, num_class, PAD_TOKEN) #model
22 loss_func = nn.CrossEntropyLoss() # loss_function
23
24 #####
25 #                               END OF YOUR CODE                               #
26 #####
27
28 # copy the model to the specified device (GPU)
29 model = model.to(device)
30
31 optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
32 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, epochs, 1e-8)
33 total_accu = None
34 train_iter, test_iter = AG_NEWS()
35 train_dataset = to_map_style_dataset(train_iter)
36 test_dataset = to_map_style_dataset(test_iter)
37 num_train = int(len(train_dataset) * 0.95)
```

```
In [13]: 1 split_train_, split_valid_ = random_split(
2         train_dataset,
3         [num_train, len(train_dataset) - num_train]
4     )
5
6 train_dataloader = DataLoader(
7     split_train_, batch_size=batch_size,
8     shuffle=True, collate_fn=collate_batch
9 )
10
11 valid_dataloader = DataLoader(
12     split_valid_, batch_size=batch_size,
13     shuffle=False, collate_fn=collate_batch
14 )
15
16 test_dataloader = DataLoader(
17     test_dataset, batch_size=batch_size,
18     shuffle=False, collate_fn=collate_batch
19 )
20 split_train_[21]
```

Out[13]: (3, 'Intel Seen Readyng New Wi-Fi Chips SAN FRANCISCO (Reuters) - Intel Corp. &A HREF="http://www.investor.reuters.com/FullQuote.aspx?ticker=INTC.0 target=/stocks/quickinfo/fullquote"&INTC.0&A&gt; this week is expected to introduce a chip that adds support for a relatively obscure version of Wi-Fi, analysts said on Monday, in a move that could help ease congestion on wireless networks.')

```
In [14]: 1 # You should be able get a validation accuracy around 86%
2 for epoch in range(1, epochs + 1):
3     # global logits_tracker
4     # logits_tracker[epoch] = None
5     epoch_start_time = time.time()
6     train(model, train_dataloader, loss_func, device, 1, optimizer)
7     accu_val = evaluate(model, valid_dataloader, loss_func, device)
8     if total_accu is not None and total_accu > accu_val:
9         scheduler.step()
10    else:
11        total_accu = accu_val
12    print('-' * 59)
13    print('| end of epoch {:3d} | time: {:.2f}s | '
14          'valid accuracy {:.3f} '.format(epoch,
15                                          time.time() - epoch_start_time,
16                                          accu_val))
17    print('-' * 59)
```

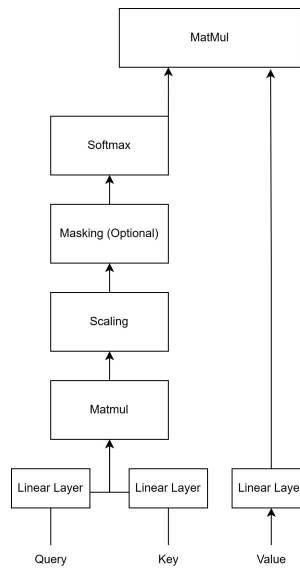
epoch	1		500/ 1782 batches		accuracy	0.559
epoch	1		1000/ 1782 batches		accuracy	0.839
epoch	1		1500/ 1782 batches		accuracy	0.876
-----						
end of epoch	1		time: 21.06s		valid accuracy	0.900
-----						
epoch	2		500/ 1782 batches		accuracy	0.907
epoch	2		1000/ 1782 batches		accuracy	0.910
epoch	2		1500/ 1782 batches		accuracy	0.910
-----						
end of epoch	2		time: 17.76s		valid accuracy	0.909
-----						
epoch	3		500/ 1782 batches		accuracy	0.929
epoch	3		1000/ 1782 batches		accuracy	0.930
epoch	3		1500/ 1782 batches		accuracy	0.928
-----						
end of epoch	3		time: 17.76s		valid accuracy	0.912
-----						

Section 2: Transformers - 'Attention is All you Need' : Classifier

Transformers are a type of deep learning architecture that has had a profound impact on a wide range of natural language processing (NLP) tasks and other sequence-to-sequence tasks. They are known for their ability to model long-range dependencies and their parallelization capabilities. A typical transformer model consists of several key components:

Output  
Probabilities

## 2.1 Multihead Attention



Multi-Head Attention can be mathematically explained as follows:

Let's assume we have a sequence of input vectors  $(X = x_1, x_2, \dots, x_n)$ , where  $(x_i)$  represents the  $(i)$ -th element of the sequence. Each  $x_i$  is typically a vector, such as a word embedding in natural language processing.

### 1. Single Attention Head:

- In a single attention head, we compute attention scores  $(A_{ij})$  between every pair of input elements  $(x_i)$  and  $(x_j)$ . These scores are computed using a compatibility function, often a dot product or a learned linear transformation followed by a softmax activation:

$$A_{ij} = \text{softmax}((Q_i x_i)^T (K_j x_j) \sqrt{d_k})$$

Where  $Q_i$  and  $K_j$  are learned linear transformations of the input vectors  $x_i$  and  $x_j$ , and  $d_k$  is the dimension of the key vectors.

- The attention scores are used to compute weighted representations of the input sequence:

$$\text{Attention}(X) = \sum_{j=1}^n A_{ij} V_j$$

Where  $V_j$  is a learned linear transformation of the input vector  $x_j$ .

### 2. Multiple Attention Heads:

- In Multi-Head Attention, we use  $H$  attention heads in parallel. Each head has its own sets of learned parameters for  $Q$ ,  $K$ , and  $V$ , resulting in  $H$  sets of attention scores and weighted representations.

$$\text{MultiHead}(X) = \text{Concatenate}(\text{Head}_1, \text{Head}_2, \dots, \text{Head}_H). W^O$$

Where  $W^O$  is another learned linear transformation applied to the concatenated outputs, and  $\text{Head}_i$  represents the output of the  $i$ -th attention head.



In [15]:

```

class MultiHeadAttention(nn.Module):
    """
    2
    A module that computes multi-head attention given query, key, and value tensors.
    """
    4
    def __init__(self, input_dim: int, num_heads: int):
        """
        6
        Constructor.
        8
        Inputs:
        10 input_dim: Dimension of the input query, key, and value. Here we assume they all have
        11 the same dimensions. But they could have different dimensions in other problems.
        12 num_heads: Number of attention heads
        13"""
        14 super(MultiHeadAttention, self).__init__()
        15
        16 assert input_dim % num_heads == 0 # Check if we can get back the original Dimensions!
        17
        18 self.input_dim = input_dim
        19 self.num_heads = num_heads
        20 self.dim_per_head = input_dim // num_heads
        21
        22 #####
        23 # TODO: Define the linear transformation layers for key, value, and query. #
        24 # Also define the output layer.
        25 #####
        26
        27 self.query_transform = nn.Linear(input_dim, input_dim) #q
        28 self.key_transform = nn.Linear(input_dim, input_dim) #k
        29 self.value_transform = nn.Linear(input_dim, input_dim) #v
        30 self.out_projection = nn.Linear(input_dim, input_dim) #out_projections
        31 #####
        32 # END OF YOUR CODE
        33 #####
        34
        35 self.scores = None
        36
    def forward(self, query: torch.Tensor, key: torch.Tensor, value: torch.Tensor, mask: torch.Tensor=None):
        """
        38
        Compute the attended feature representations.
        40
        Inputs:
        42 query: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
        43 and C is the channel dimension
        44 key: Tensor of the shape BxLxC
        45 value: Tensor of the shape BxLxC
        46 mask: Tensor indicating where the attention should *not* be performed
        47"""
        48 b = query.shape[0]
        49
        50 dot_prod_scores = None
        51 #####
        52 # TODO: Compute the scores based on dot product between transformed query, #
        53 # key, and value. You may find torch.matmul helpful, whose documentation #
        54 # can be found at
        55 # https://pytorch.org/docs/stable/generated/torch.matmul.html#torch.matmul#
        56 # Remember to divide the dot product similarity scores by square root of #
        57 # the channel dimension per head.
        58 #
        59 # Since no for loops are allowed here, think of how to use tensor reshape #
        60 # to process multiple attention heads at the same time.
        61 #####
        62
        63 # Transform query, key, and value
        64 query_trans = self.query_transform(query) # transform q
        65 key_trans = self.key_transform(key) # transform k
        66 value_trans = self.value_transform(value) # transform v
        67
        68 # Reshape to allow multiple heads
        69 query_trans = query_trans.view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
        70 key_trans = key_trans.view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
        71 value_trans = value_trans.view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
        72
        73
        74 dot_prod_scores = torch.matmul(query_trans, key_trans.transpose(-2, -1)) # Computing dot product attention scores
        75
        76
        77 dot_prod_scores = dot_prod_scores / math.sqrt(self.dim_per_head) #scaling
        78
        79
        80 #####
        81 # END OF YOUR CODE
        82 #####
        83
        84 if mask is not None:
        85     # We simply set the similarity scores to be near zero for the positions
        86     # where the attention should not be done. Think of why we do this.
        87     dot_prod_scores = dot_prod_scores.masked_fill(mask == 0, -1e9) #, batch, h, LxL
        88
        89 out = None
        90 #####
        91 # TODO: Compute the attention scores, which are then used to modulate the #
        92 # value tensor. Finally concatenate the attended tensors from multiple heads #
        93 # and feed it into the output layer. You may still find torch.matmul #
        94 # helpful.
        95 #
        96 # Again, think of how to use reshaping tensor to do the concatenation. #
        97 #####
        98
        99 raise NotImplementedError
        100 attention_weights = torch.softmax(dot_prod_scores, dim=-1) # BxHxLxL
        101
        102 # Compute attended values
        103 out = torch.matmul(attention_weights, value_trans) # BxHxLxC'
        104
        105
        106 # Reshape and concatenate the heads
        107 out = out.transpose(1, 2).contiguous().view(b, -1, self.input_dim) # BxLxHxC' -> BxLxC
        108 out = self.out_projection(out) # Apply output projection
        109 #####
        110 # END OF YOUR CODE
        111 #####
        112
        113 return out

```

```
In [16]: 1 # Sanity Check
2 x = torch.randn((2, 10, 8))
3 mask = torch.randn((2, 10)) > 0.5
4 mask = mask.unsqueeze(1).unsqueeze(-1)
5 num_heads = 4
6 model = MultiHeadAttention(8, num_heads)
7 y = model(x, x, x, mask)
8 assert len(y.shape) == len(x.shape)
9 for dim_x, dim_y in zip(x.shape, y.shape):
10     assert dim_x == dim_y
```

2.2 Positional Encoding Module

Positional Encoding is a critical component in the Transformer architecture, designed to provide information about the positions of elements in a sequence to a model that inherently lacks sequential information. Transformers use self-attention mechanisms that do not inherently understand the order or position of tokens in the input. Positional Encoding is introduced to address this limitation and allow the model to consider the order of elements within the input sequence.

It addresses the challenge of modeling sequences with self-attention mechanisms that do not inherently understand the order of elements. By adding Positional Encoding to the input embeddings, the model can differentiate between tokens based on their positions and capture sequential information effectively.

1. Positional Encoding Function:

- Positional Encoding is typically represented as a fixed-size vector that is added element-wise to the input embeddings. This vector is determined by a mathematical function.
- The most common approach is to use a combination of sine and cosine functions with different frequencies and phases to create a unique encoding for each position.
- For each position *pos* and dimension *i* of the Positional Encoding vector, *PE(pos, 2i)* is given by

$$\sin\left(\frac{pos}{10000^{(2i/d_{model})}}\right)$$
 if *i* is even  
$$\cos\left(\frac{pos}{10000^{(2i/d_{model})}}\right)$$
 if *i* is odd

*d<sub>model</sub>* is the dimension of the model's input embeddings.

2. Adding Positional Encoding:

- The Positional Encoding vector is added element-wise to the input embeddings. This combination of the original word embeddings and the Positional Encoding allows the model to distinguish between tokens based on their positions.

2.2.1 Let's Try to work an example!

Assume the sentence coming into the encoding is, "This is an Example"; The Positional Encoding layer is initialized with the following parameters;

- *k* = 0 ≤ *k* < *L* = 3
- *n* = 100
- *d* = 6
- *I* = 0 ≤ *i* < *d*/2 = 2
- Based on the values above, for the text given, Find the Position Encoding values below or create on on your own and upload to this section! **DONT CODE IT**

d= 6 & n = 100

Sequence	Index(k)	i=0	i=0	i=1	i=1	i=2	i=2	i=3	i=3
This	0	P(00) = 0	1	0	1	0	1	0	1
Is	1	P(10) = 0.841	0.540	0.213	0.976	0.046	0.998	0.009	0.999
An	2	P(20) = 0.909	-0.416	0.417	0.908	0.092	0.995	0.019	0.999
Example	3	P(30) = 0.141	-0.990	0.602	0.798	0.138	0.990	0.029	0.999

2.2.1 Let's Code!

Now try to use the same approach to implement the Positional Encoding part.

For full credit do not use for loops;

Make use of packages like

- torch.arange() : <https://pytorch.org/docs/stable/generated/torch.arange.html> (<https://pytorch.org/docs/stable/generated/torch.arange.html>)
- torch.stack() : <https://pytorch.org/docs/stable/generated/torch.stack.html> (<https://pytorch.org/docs/stable/generated/torch.stack.html>)

```

In [17]: 1 class PositionalEncoding(nn.Module):
2         """
3         A module that adds positional encoding to each of the token's features.
4         So that the Transformer is position aware.
5         """
6         def __init__(self, input_dim: int, max_len: int=10000):
7             """
8             Inputs:
9             - input_dim: Input dimension about the features for each token
10            - max_len: The maximum sequence length
11            """
12            super(PositionalEncoding, self).__init__()
13
14            self.input_dim = input_dim
15            self.max_len = max_len
16
17
18        def forward(self, x):
19            """
20            Compute the positional encoding and add it to x.
21
22            Input:
23            - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
24              and C is the channel dimension
25
26            Return:
27            - x: Tensor of the shape BxLxC, with the positional encoding added to the input
28            """
29            seq_len = x.shape[1]
30            input_dim = x.shape[2]
31
32            #####
33            # TODO: Compute the positional encoding
34            # Check Section 3.5 for the definition (https://arxiv.org/pdf/1706.03762.pdf)
35            #
36            # It's a bit messy, but the definition is provided for your here for your
37            # convenience (in LaTeX).
38            #  $PE_{\{pos, 2i\}} = \sin(pos / 10000^{2i/d_{model}})$ 
39            #  $PE_{\{pos, 2i+1\}} = \cos(pos / 10000^{2i/d_{model}})$ 
40            #
41            # You should replace 10000 with max_len here.
42            #####
43
44            pos_indices = torch.arange(seq_len).unsqueeze(1)
45            divisors = torch.pow(self.max_len, torch.arange(0, input_dim, 2).float() / input_dim)
46            pe = torch.zeros(seq_len, input_dim)
47            pe[:, 0::2] = torch.sin(pos_indices.float() / divisors)
48            pe[:, 1::2] = torch.cos(pos_indices.float() / divisors)
49
50
51            #####
52            #                               END OF YOUR CODE
53            #####
54
55            x = x + pe.to(x.device)
56            return x

```

```

In [18]: 1 # Sanity check - I
2 x = torch.randn(1, 100, 20)
3 pe = PositionalEncoding(20)
4 y = pe(x)
5 print(y)
6 assert len(x.shape) == len(y.shape)
7 for dim_x, dim_y in zip(x.shape, y.shape):
8     assert dim_x == dim_y

```

```

tensor([[[[-0.1095, -0.5085, -0.4006, ..., 0.2879, -0.9944, 2.2845],
[ 1.9422, -0.3290, 0.6247, ..., 2.3464, 0.5444, 1.6659],
[ 0.9008, -1.3341, 1.7986, ..., 0.2932, -0.3308, 0.8616],
...,
[-0.5053, -1.2110, 0.4801, ..., 2.3540, 0.3029, 0.6670],
[-1.0749, -1.8689, 0.4923, ..., 1.5702, -0.2510, -0.1182],
[-0.2595, 0.6789, 0.9136, ..., 2.1194, -0.3044, 0.9322]]]])

```

```

In [19]: 1 # Sanity Check - II
2 x = torch.randn(1, 100, 6)
3 d = 6
4 n = 100
5 pe = PositionalEncoding(d,n)
6 y = pe(x)
7
8 y -= x
9 print(y[:, :4, :])

```

```

tensor([[[[ 0.0000, 1.0000, 0.0000, 1.0000, 0.0000, 1.0000],
[ 0.8415, 0.5403, 0.2138, 0.9769, 0.0464, 0.9989],
[ 0.9093, -0.4161, 0.4177, 0.9086, 0.0927, 0.9957],
[ 0.1411, -0.9900, 0.6023, 0.7983, 0.1388, 0.9903]]]])

```

```

tensor([[[[ 0.0000, 1.0000, 0.0000, 1.0000, 0.0000, 1.0000],
[ 0.8415, 0.5403, 0.2138, 0.9769, 0.0464, 0.9989],
[ 0.9093, -0.4161, 0.4177, 0.9086, 0.0927, 0.9957],
[ 0.1411, -0.9900, 0.6023, 0.7983, 0.1388, 0.9903]]]])

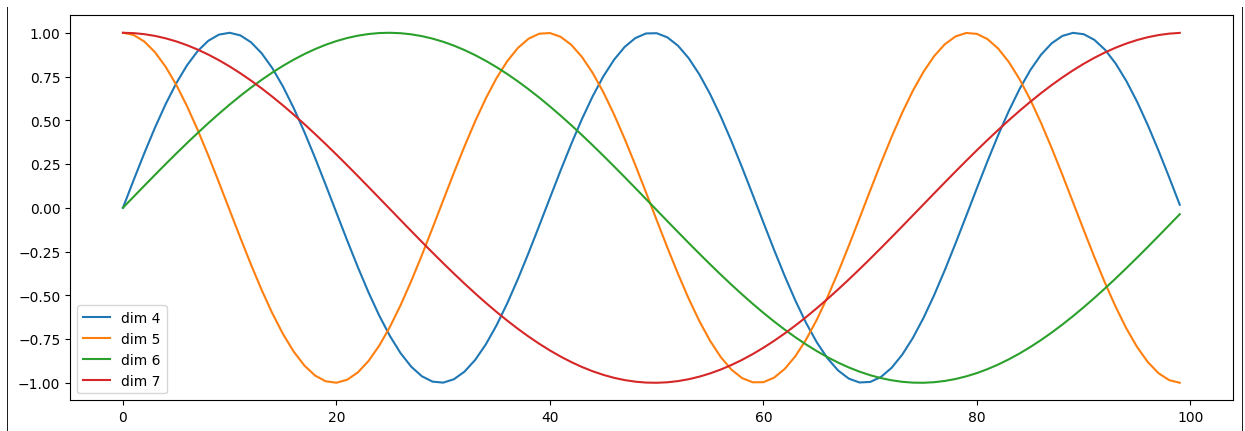
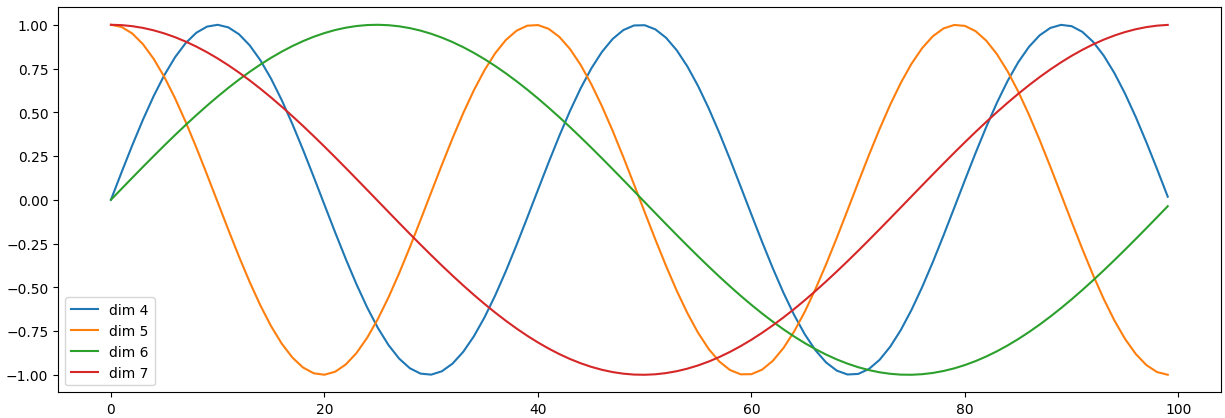
```

```

In [20]: 1 # Sanity check - III
          2 %matplotlib inline
          3 import matplotlib.pyplot as plt
          4 import numpy as np
          5
          6 plt.figure(figsize=(15, 5))
          7 pe = PositionalEncoding(20)
          8 y = pe.forward(torch.zeros(1, 100, 20))
          9 plt.plot(np.arange(100), y[0, :, 4:8].data.numpy())
         10 plt.legend(["dim %d"%p for p in [4,5,6,7]])

```

Out[20]: <matplotlib.legend.Legend at 0x7da92c9bbe80>



## 2.3 FeedForward Module

The FeedForward Layer in the Transformer architecture is a position-wise neural network layer designed to process the context-aware representations generated by the self-attention mechanism. It consists of two linear transformations followed by a non-linear activation function, typically ReLU. The FeedForward Layer is applied independently to each position in the sequence, allowing the model to capture different patterns at different positions. This position-wise independence, combined with non-linearity, helps the model learn complex relationships within the data and plays a crucial role in the Transformer's ability to process and understand sequential data effectively, making it a fundamental component for various sequence-to-sequence tasks.

Mathematically, if  $X$  represents the input sequence (a sequence of embeddings),  $FFN(X)$  is the output of the FeedForward Layer, and  $W_1$ ,  $W_2$ ,  $b_1$ , and  $b_2$  represent learned weight matrices and bias terms, the operation can be expressed as,

$$FFN(X) = \text{ReLU}(X \cdot W_1 + b_1) \cdot W_2 + b_2.$$

```

In [21]: 1 class FeedForwardNetwork(nn.Module):
2         """
3         A simple feedforward network. Essentially, it is a two-layer fully-connected
4         neural network.
5         """
6         def __init__(self, input_dim, ff_dim):
7             """
8             Inputs:
9             - input_dim: Input dimension
10            - ff_dim: Hidden dimension
11            """
12            super(FeedForwardNetwork, self).__init__()
13
14            #####
15            # TODO: Define the two linear layers and a non-linear one.
16            #####
17
18            self.linear1 = nn.Linear(input_dim, ff_dim)
19            self.relu = nn.ReLU()
20            self.linear2 = nn.Linear(ff_dim, input_dim)
21
22            #####
23            #                               END OF YOUR CODE                               #
24            #####
25
26        def forward(self, x: torch.Tensor):
27            """
28            Input:
29            - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
30              and C is the channel dimension
31
32            Return:
33            - y: Tensor of the shape BxLxC
34            """
35
36            #y = None
37            #####
38            # TODO: Process the input.
39            #####
40
41            y = self.linear1(x)
42            y = self.relu(y)
43            y = self.linear2(y)
44            #####
45            #                               END OF YOUR CODE                               #
46            #####
47
48            return y
49

```

```

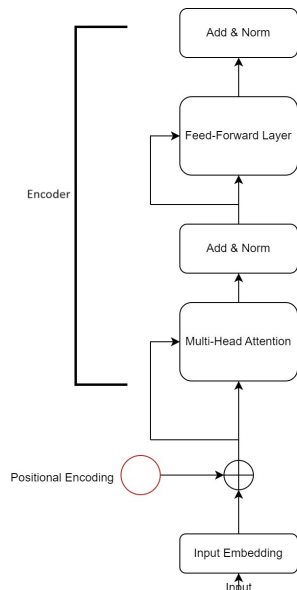
In [22]: 1 # Sanity Check
2 x = torch.randn((2, 10, 8))
3 ff_dim = 4
4 model = FeedForwardNetwork(8, ff_dim)
5 y = model(x)
6 assert len(x.shape) == len(y.shape)
7 for dim_x, dim_y in zip(x.shape, y.shape):
8     assert dim_x == dim_y
9 print(y.shape)

torch.Size([2, 10, 8])

```

## 2.4 Encoder Module

The Encoder module in a Transformer is responsible for processing the input sequence, typically used for tasks like language understanding and representation learning. It consists of multiple identical layers, each containing two main components: the Multi-Head Self-Attention mechanism and the Position-wise FeedForward Layer.



- In each layer, the input sequence is first passed through the Multi-Head Self-Attention mechanism, which computes weighted representations for each element in the sequence, capturing contextual information. The attention output is then passed through the Position-wise FeedForward Layer, introducing non-linearity and allowing the model to capture different patterns at each position.
- This process is repeated for each layer in the encoder stack, enabling the model to capture hierarchical features and within the input sequence effectively. The final encoder output represents a rich contextualized representation of the input sequence, which can be used for various downstream tasks, including translation, text generation, and sentiment analysis.

## 2.4.1 Encoder Cell

```

In [23]: 1 class TransformerEncoderCell(nn.Module):
2         """
3         A single cell (unit) for the Transformer encoder.
4         """
5         def __init__(self, input_dim: int, num_heads: int, ff_dim: int, dropout: float):
6             """
7             Inputs:
8             - input_dim: Input dimension for each token in a sequence
9             - num_heads: Number of attention heads in a multi-head attention module
10            - ff_dim: The hidden dimension for a feedforward network
11            - dropout: Dropout ratio for the output of the multi-head attention and feedforward
12                      modules.
13            """
14            super(TransformerEncoderCell, self).__init__()
15
16            #####
17            # TODO: A single Transformer encoder cell consists of
18            # 1. A multi-head attention module
19            # 2. Followed by dropout
20            # 3. Followed by layer norm (check nn.LayerNorm)
21            # https://pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html#torch.nn.LayerNorm
22            #
23            # At the same time, it also has
24            # 1. A feedforward network
25            # 2. Followed by dropout
26            # 3. Followed by layer norm
27            #####
28
29            self.multi_head_attn = MultiHeadAttention(input_dim, num_heads)
30            self.attn_dropout = nn.Dropout(dropout)
31            self.attn_layer_norm = nn.LayerNorm(input_dim)
32            self.identity = nn.Identity()
33
34            self.ffn = FeedForwardNetwork(input_dim, ff_dim)
35
36            self.ffn_dropout = nn.Dropout(dropout)
37            self.ffn_layer_norm = nn.LayerNorm(input_dim)
38            #####
39            # END OF YOUR CODE
40            #####
41            self.attention = None
42
43            def forward(self, x: torch.Tensor, mask: torch.Tensor=None):
44                """
45                Inputs:
46                - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
47                  and C is the channel dimension
48                - mask: Tensor for multi-head attention
49                """
50
51                y = None
52                #####
53                # TODO: Get the output of the multi-head attention part (with dropout
54                # and layer norm), which is used as input to the feedforward network (
55                # again, followed by dropout and layer norm).
56                #
57                # Don't forget the residual connections for both parts. Append the
58                # 1st Normalized Output before feed_forward to self.attention(Useful in
59                # visualizing)
60                #####
61
62                # multihead attention
63                skip = self.identity(x)
64                x = self.multi_head_attn(x, x, x, mask)
65                x = x + skip
66                x = self.attn_dropout(x)
67                self.attention = self.attn_layer_norm(x)
68
69                # feed forward network
70                skip = self.identity(self.attention)
71                x = self.ffn(self.attention)
72                x = x + skip
73                x = self.ffn_dropout(x)
74                y = self.ffn_layer_norm(x)
75
76                #####
77                # END OF YOUR CODE
78                #####
79
80            return y

```

```

In [24]: 1 # Sanity Check
2
3 x = torch.randn((2, 10, 8))
4 mask = torch.randn((2, 10)) > 0.5
5 mask = mask.unsqueeze(1).unsqueeze(-1)
6 num_heads = 4
7 model = TransformerEncoderCell(8, num_heads, 32, 0.1)
8 y = model(x, mask)
9 assert len(x.shape) == len(y.shape)
10 for dim_x, dim_y in zip(x.shape, y.shape):
11     assert dim_x == dim_y
12 print(y.shape)

```

torch.Size([2, 10, 8])

## 2.4.2 Building an Encoder Module

```
In [25]: 1 class TransformerEncoder(nn.Module):
2         """
3         A full encoder consisting of a set of TransformerEncoderCell.
4         """
5         def __init__(self, input_dim: int, num_heads: int, ff_dim: int, num_cells: int, dropout: float=0.1):
6             """
7             Inputs:
8             - input_dim: Input dimension for each token in a sequence
9             - num_heads: Number of attention heads in a multi-head attention module
10            - ff_dim: The hidden dimension for a feedforward network
11            - num_cells: Number of TransformerEncoderCells
12            - dropout: Dropout ratio for the output of the multi-head attention and feedforward
13              modules.
14            """
15            super(TransformerEncoder, self).__init__()
16
17            #####
18            # TODO: Construct a nn.ModuleList to store a stack of      #
19            # TransformerEncoderCells. Check the documentation here of how to use it #
20            # https://pytorch.org/docs/stable/generated/torch.nn.ModuleList.html#torch.nn.ModuleList
21
22            # At the same time, define a layer normalization layer to process the      #
23            # output of the entire encoder.                                           #
24            #####
25
26            self.norm = nn.LayerNorm(input_dim)
27            self.cells = nn.ModuleList([TransformerEncoderCell(input_dim, num_heads, ff_dim, dropout) for _ in range(num_cells)])
28
29            #####
30            #                               END OF YOUR CODE                               #
31            #####
32
33            def forward(self, x: torch.Tensor, mask: torch.Tensor=None):
34                """
35                Inputs:
36                - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
37                  and C is the channel dimension
38                - mask: Tensor for multi-head attention
39
40                Return:
41                - y: Tensor of the shape of BxLxC, which is the normalized output of the encoder
42                """
43
44                #####
45                # TODO: Feed x into the stack of TransformerEncoderCells and then      #
46                # normalize the output with layer norm.                               #
47                #####
48                for cell in self.cells:
49                    x = cell(x, mask)
50                y = self.norm(x)
51                #####
52                #                               END OF YOUR CODE                               #
53                #####
54
55                return y
56
```

## 2.5 Transformer Classifier

Now, lets put this all the above described modules together to make out classifier

```

In [26]: 1 class TransformerClassifier(nn.Module):
          2     """
          3     A Transformer-based text classifier.
          4     """
          5     def __init__(self,
          6                 vocab_size: int, embed_dim: int, num_heads: int, trx_ff_dim: int,
          7                 num_trx_cells: int, num_class: int, dropout: float=0.1, pad_token: int=0
          8                 ):
          9         """
         10         Inputs:
         11         - vocab_size: Vocabulary size, indicating how many tokens we have in total.
         12         - embed_dim: The dimension of word embeddings
         13         - num_heads: Number of attention heads in a multi-head attention module
         14         - trx_ff_dim: The hidden dimension for a feedforward network
         15         - num_trx_cells: Number of TransformerEncoderCells
         16         - dropout: Dropout ratio
         17         - pad_token: The index of the padding token.
         18         """
         19         super(TransformerClassifier, self).__init__()
         20
         21         self.embed_dim = embed_dim
         22
         23         # word embedding layer
         24         self.embedding = nn.Embedding(vocab_size, embed_dim, padding_idx=pad_token)
         25
         26         #####
         27         # TODO: Define a module for positional encoding, Transformer encoder, and #
         28         # a output layer
         29         #####
         30
         31         self.positional_encoding = PositionalEncoding(embed_dim)
         32
         33         self.encoder = TransformerEncoder(embed_dim, num_heads, trx_ff_dim, num_trx_cells, dropout)
         34
         35         self.fc = nn.Linear(embed_dim, num_class)
         36         #####
         37         #                               END OF YOUR CODE
         38         #####
         39
         40     def forward(self, text, mask=None):
         41         """
         42         Inputs:
         43         - text: Tensor with the shape of BxLxC.
         44         - mask: Tensor for multi-head attention
         45
         46         Return:
         47         - logits: Tensor with the shape of BxK, where K is the number of classes
         48         """
         49
         50         # word embeddings, note we multiple the embeddings by a factor
         51         embedded = self.embedding(text) * math.sqrt(self.embed_dim)
         52
         53         #####
         54         # TODO: Apply positional embedding to the input, which is then fed into #
         55         # the encoder. Average pooling is applied then to all the features of all #
         56         # tokens. Finally, the logits are computed based on the pooled features. #
         57         #####
         58
         59         embedded = self.positional_encoding(embedded)
         60
         61         encoder_output = self.encoder(embedded, mask)
         62
         63         pooled = encoder_output.mean(dim=1)
         64
         65         logits = self.fc(pooled)
         66         #####
         67         #                               END OF YOUR CODE
         68         #####
         69
         70         return logits

```

```

In [27]: 1 # Sanity Check
          2 vocab_size = 10
          3 embed_dim = 16
          4 num_heads = 4
          5 trx_ff_dim = 16
          6 num_trx_cells = 2
          7 num_class = 3
          8
          9 x = torch.arange(vocab_size).view(1, -1)
         10 x = torch.cat((x, x), dim=0)
         11
         12 mask = (x != 0).unsqueeze(-2).unsqueeze(1)
         13 model = TransformerClassifier(vocab_size, embed_dim, num_heads, trx_ff_dim, num_trx_cells, num_class)
         14 print('x: {}, mask: {}'.format(x.shape, mask.shape))
         15 y = model(x, mask)
         16 assert len(y.shape) == 2 and y.shape[0] == x.shape[0] and y.shape[1] == num_class
         17 print(y.shape)

```

x: torch.Size([2, 10]), mask: torch.Size([2, 1, 1, 10])  
torch.Size([2, 3])

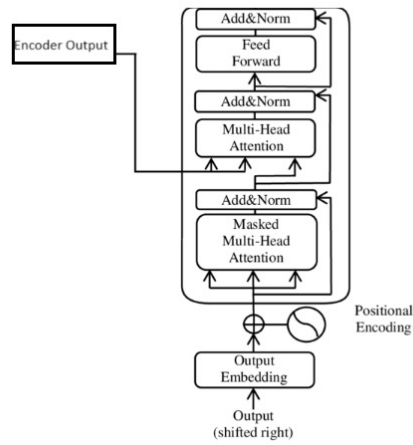


## 2.6 Deciding HyperParameters & Training

```
In [28]: 1 assert torch.cuda.is_available()
2
3 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
4
5 # Hyperparameters
6 epochs = 3 # epoch
7 lr = 0.0005 # learning rate
8 batch_size = 64 # batch size for training
9
10 train_iter = AG_NEWS(split='train')
11 num_class = len(set([label for (label, text) in train_iter]))
12 vocab_size = len(vocab)
13 emsize = 64
14
15 num_heads = 4
16 num_trx_cells = 2
17
18 gradient_norm_clip = 1
19
20 #####
21 # Define a Transformer-based text classifier and a loss function. #
22 #####
23
24 model = TransformerClassifier(
25     vocab_size=vocab_size,
26     embed_dim=emsize,
27     num_heads=num_heads,
28     trx_ff_dim=emsize,
29     num_trx_cells=num_trx_cells,
30     num_class=num_class,
31     dropout=0.1,
32     pad_token=vocab['<pad>']
33 )
34
35 loss_func = nn.CrossEntropyLoss()
36 #####
37 #                               END OF YOUR CODE                               #
38 #####
39 model = model.to(device)
40
41 optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
42 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, epochs, 1e-8)
43 total_accu = None
44
45 # You should be able to get a validation accuracy around 89%
46 for epoch in range(1, epochs + 1):
47     epoch_start_time = time.time()
48     train(model, train_dataloader, loss_func, device, gradient_norm_clip, optimizer)
49     accu_val = evaluate(model, valid_dataloader, loss_func, device)
50     if total_accu is not None and total_accu > accu_val:
51         scheduler.step()
52     else:
53         total_accu = accu_val
54     print('-' * 59)
55     print('| end of epoch {:3d} | time: {:.2f}s | '
56           'valid accuracy {:.3f}'.format(epoch,
57                                         time.time() - epoch_start_time,
58                                         accu_val))
59     print('-' * 59)
```

epoch	1		500/ 1782 batches		accuracy	0.651
epoch	1		1000/ 1782 batches		accuracy	0.784
epoch	1		1500/ 1782 batches		accuracy	0.826
-----						
end of epoch	1		time: 29.22s		valid accuracy	0.851
-----						
epoch	2		500/ 1782 batches		accuracy	0.864
epoch	2		1000/ 1782 batches		accuracy	0.871
epoch	2		1500/ 1782 batches		accuracy	0.879
-----						
end of epoch	2		time: 27.62s		valid accuracy	0.882
-----						
epoch	3		500/ 1782 batches		accuracy	0.895
epoch	3		1000/ 1782 batches		accuracy	0.901
epoch	3		1500/ 1782 batches		accuracy	0.902
-----						
end of epoch	3		time: 27.78s		valid accuracy	0.888

## Section 3: Transformers - 'Attention is All you need' : Machine Translation



The decoder module in a Seq-Seg Transformer model is responsible for generating the output sequence based on the information gathered by the encoder and previous tokens in an autoregressive manner. It utilizes self-attention mechanisms with masking to enforce causality and multi-head attention to capture dependencies between tokens in the output. Self-attention calculates attention weights for each position in the sequence, and multi-head attention aggregates the results from multiple attention heads, enhancing the model's representational power. Cross-attention is also employed to allow the decoder to focus on relevant parts of the encoder's output.

Additionally, position-wise feed-forward networks further process the information by applying linear transformations and non-linear activation functions to each position independently. Throughout the decoder, layer normalization and residual connections are utilized to enhance training stability. These components are typically stacked in multiple layers to enable the model to learn complex relationships and generate coherent output sequences.

## 3.1.1 Decoder Cell

```

In [29]: 1 class TransformerDecoderCell(nn.Module):
2         """
3         A single cell (unit) of the Transformer decoder.
4         """
5         def __init__(self, input_dim: int, num_heads: int, ff_dim: int, dropout: float=0.1):
6             """
7             Inputs:
8             - input_dim: Input dimension for each token in a sequence
9             - num_heads: Number of attention heads in a multi-head attention module
10            - ff_dim: The hidden dimension for a feedforward network
11            - dropout: Dropout ratio for the output of the multi-head attention and feedforward
12              modules.
13            """
14            super(TransformerDecoderCell, self).__init__()
15
16            #####
17            # TODO: Similar to the TransformerEncoderCell, define two #
18            # MultiHeadAttention modules. One for processing the tokens on the #
19            # decoder side. The other for getting the attention across the encoder. #
20            # and the decoder. Also define a feedforward network. Don't forget the #
21            # Dropout and Layer Norm layers. #
22            #####
23
24            self.self_attn = MultiHeadAttention(input_dim, num_heads)
25            self.self_attn_dropout = nn.Dropout(dropout)
26            self.self_attn_layer_norm = nn.LayerNorm(input_dim)
27            self.cross_attn = MultiHeadAttention(input_dim, num_heads)
28            self.cross_attn_dropout = nn.Dropout(dropout)
29            self.cross_attn_layer_norm = nn.LayerNorm(input_dim)
30            self.identity = nn.Identity()
31            self.ffn = FeedForwardNetwork(input_dim, ff_dim)
32            self.ffn_dropout = nn.Dropout(dropout)
33            self.ffn_layer_norm = nn.LayerNorm(input_dim)
34
35            #####
36            #                               END OF YOUR CODE                               #
37            #####
38
39            def forward(self, x: torch.Tensor, encoder_output: torch.Tensor, src_mask=None, tgt_mask=None):
40                """
41                Inputs:
42                - x: Tensor of BxLxC, word embeddings on the decoder side
43                - encoder_output: Tensor of BxLxC, word embeddings on the encoder side
44                - src_mask: Tensor, masks of the tokens on the encoder side
45                - tgt_mask: Tensor, masks of the tokens on the decoder side
46
47                Return:
48                - y: Tensor of BxLxC. Attended features for all tokens on the decoder side.
49                """
50
51                y = None
52                #####
53                # TODO: Compute the self-attended features for the tokens on the decoder #
54                # side. Then compute the cross-attended features for the tokens on the #
55                # decoder side to the encoded features, which are finally feed into the #
56                # feedforward network #
57                #####
58
59                # Self-attention
60                skip = self.identity(x)
61                x = self.self_attn(x, x, x, tgt_mask)
62                x = self.self_attn_dropout(x)
63                x = x + skip
64                x = self.self_attn_layer_norm(x)
65
66                # cross-attention
67                skip = self.identity(x)
68                x = self.cross_attn(x, encoder_output, encoder_output, src_mask)
69                x = self.cross_attn_dropout(x)
70                x = x + skip
71                x = self.cross_attn_layer_norm(x)
72
73                # Feed Forward Networks
74                skip = self.identity(x)
75                x = self.ffn(x)
76                x = self.ffn_dropout(x)
77                x = x + skip
78                y = self.ffn_layer_norm(x)
79
80                #####
81                #                               END OF YOUR CODE                               #
82                #####
83
84                return y
85

```

```

In [30]: 1 # Sanity Check
2
3 dec_feats = torch.randn((3, 10, 16))
4 dec_mask = torch.randn((3, 1, 10, 10)) > 0.5
5
6 enc_feats = torch.randn((3, 12, 16))
7 enc_mask = torch.randn((3, 1, 1, 12)) > 0.5
8
9 model = TransformerDecoderCell(16, 2, 32, 0.1)
10
11 z = model(dec_feats, enc_feats, enc_mask, dec_mask)
12 assert len(z.shape) == len(dec_feats.shape)
13 for dim_z, dim_x in zip(z.shape, dec_feats.shape):
14     assert dim_z == dim_x
15 print(z.shape)

```

```
torch.Size([3, 10, 16])
```

## 3.1.2 Building the Decoder Module

```

In [31]: 1 class TransformerDecoder(nn.Module):
          2     """
          3     A TransformerDecoder is a stack of multiple TransformerDecoderCells and a Layer Norm.
          4     """
          5     def __init__(self, input_dim: int, num_heads: int, ff_dim: int, num_cells: int, dropout=0.1):
          6         """
          7         Inputs:
          8         - input_dim: Input dimension for each token in a sequence
          9         - num_heads: Number of attention heads in a multi-head attention module
         10         - ff_dim: The hidden dimension for a feedforward network
         11         - num_cells: How many TransformerDecoderCells in stack
         12         - dropout: Dropout ratio for the output of the multi-head attention and feedforward
         13           modules.
         14         """
         15         super(TransformerDecoder, self).__init__()
         16
         17         #####
         18         # TODO: Construct a nn.ModuleList to store a stack of
         19         # TransformerDecoderCells. Check the documentation here of how to use it
         20         # https://pytorch.org/docs/stable/generated/torch.nn.ModuleList.html#torch.nn.ModuleList
         21
         22         # At the same time, define a layer normalization layer to process the
         23         # output of the entire encoder.
         24         #####
         25
         26         self.norm = nn.LayerNorm(input_dim)
         27
         28         self.cells = nn.ModuleList([TransformerDecoderCell(input_dim, num_heads, ff_dim, dropout) for _ in range(num_cells)])
         29
         30         #####
         31         #                               END OF YOUR CODE
         32         #####
         33
         34     def forward(self, x: torch.Tensor, encoder_output: torch.Tensor, src_mask=None, tgt_mask=None):
         35         """
         36         Inputs:
         37         - x: Tensor of BxLdxC, word embeddings on the decoder side
         38         - encoder_output: Tensor of BxLexC, word embeddings on the encoder side
         39         - src_mask: Tensor, masks of the tokens on the encoder side
         40         - tgt_mask: Tensor, masks of the tokens on the decoder side
         41
         42         Return:
         43         - y: Tensor of BxLdxC. Attended features for all tokens on the decoder side.
         44         """
         45
         46         y = None
         47         #####
         48         # TODO: Feed x into the stack of TransformerDecoderCells and then
         49         # normalize the output with layer norm.
         50         #####
         51
         52         for cell in self.cells:
         53             x = cell(x, encoder_output, src_mask, tgt_mask)
         54
         55         y = self.norm(x)
         56         #####
         57         #                               END OF YOUR CODE
         58         #####
         59
         60         return y

```

```

In [32]: 1 # Sanity Check
          2 dec_feats = torch.randn((3, 10, 16))
          3 dec_mask = torch.randn((3, 1, 10, 10)) > 0.5
          4
          5 enc_feats = torch.randn((3, 12, 16))
          6 enc_mask = torch.randn((3, 1, 1, 12)) > 0.5
          7
          8 model = TransformerDecoder(16, 2, 32, 2, 0.1)
          9 z = model(dec_feats, enc_feats, enc_mask, dec_mask)
         10 assert len(z.shape) == len(dec_feats.shape)
         11 for dim_z, dim_x in zip(z.shape, dec_feats.shape):
         12     assert dim_z == dim_x
         13 print(z.shape)

```

torch.Size([3, 10, 16])



### 3.2 Transformer Based Seq-to-Seq Model

```
In [33]: 1 class Seq2SeqTransformer(nn.Module):
2         """
3         Transformer-based sequence-to-sequence model.
4         """
5         def __init__(self,
6             num_encoder_layers: int, num_decoder_layers: int, embed_dim: int,
7             num_heads: int, src_vocab_size: int, tgt_vocab_size: int,
8             trx_ff_dim: int = 512, dropout: float = 0.1, pad_token: int=0
9         ):
10            """
11            Inputs:
12            - num_encoder_layers: How many TransformerEncoderCell in stack
13            - num_decoder_layers: How many TransformerDecoderCell in stack
14            - embed_dim: Word embeddings dimension
15            - num_heads: Number of attention heads
16            - src_vocab_size: Number of tokens in the source language vocabulary
17            - tgt_vocab_size: Number of tokens in the target language vocabulary
18            - trx_ff_dim: Hidden dimension in the feedforward network
19            - dropout: Dropout ratio
20            """
21            super(Seq2SeqTransformer, self).__init__()
22
23            self.embed_dim = embed_dim
24
25            # Word embeddings for both the source and target languages
26            self.src_token_embed = nn.Embedding(src_vocab_size, embed_dim, padding_idx=pad_token)
27            self.tgt_token_embed = nn.Embedding(tgt_vocab_size, embed_dim, padding_idx=pad_token)
28
29            #####
30            # TODO: Define the positional encoding, encoder, decoder, and the output #
31            # layer. Think of how many classes are in the output layer.           #
32            #####
33
34            self.positional_encoding = PositionalEncoding(embed_dim)
35
36            self.transformer_encoder = TransformerEncoder(embed_dim, num_heads, trx_ff_dim, num_encoder_layers, dropout)
37
38            self.transformer_decoder = TransformerDecoder(embed_dim, num_heads, trx_ff_dim, num_decoder_layers, dropout)
39
40            self.output_layer = nn.Linear(embed_dim, tgt_vocab_size)
41
42            #####
43            #                               END OF YOUR CODE                               #
44            #####
45
46            def forward(self, src: torch.Tensor, tgt: torch.Tensor, src_mask: torch.Tensor, tgt_mask: torch.Tensor):
47                """
48                Inputs:
49                - src: Tensor of BxLe, word indexes in the source language
50                - tgt: Tensor of BxLd, word indexes in the target language
51                - src_mask: Tensor, masks of the tokens on the encoder side
52                - tgt_mask: Tensor, masks of the tokens on the decoder side
53
54                Return:
55                - y: Tensor of BxLdxK. K is the number of classes in the output.
56                """
57
58                # Get word embeddings. Note they are scaled.
59                src_embed = self.src_token_embed(src) * math.sqrt(self.embed_dim)
60                tgt_embed = self.tgt_token_embed(tgt) * math.sqrt(self.embed_dim)
61
62                logits = None
63                #####
64                # TODO: Add positional encodings to the word embeddings. Feed them then #
65                # to the encoder and decoder, respectively. Get the logits finally.       #
66                #####
67
68                src_embed = self.positional_encoding(src_embed)
69                tgt_embed = self.positional_encoding(tgt_embed)
70
71                memory = self.transformer_encoder(src_embed, src_mask)
72                out = self.transformer_decoder(tgt_embed, memory, src_mask, tgt_mask)
73
74                logits = self.output_layer(out)
75                #####
76                #                               END OF YOUR CODE                               #
77                #####
78
79                return logits
80
81            def encode(self, src: torch.Tensor, src_mask: torch.Tensor):
82                src_embed = self.src_token_embed(src) * math.sqrt(self.embed_dim)
83                return self.transformer_encoder(self.positional_encoding(src_embed), src_mask)
84
85            def decode(self, tgt: torch.Tensor, memory: torch.Tensor, src_mask: torch.Tensor, tgt_mask: torch.Tensor):
86                tgt_embed = self.tgt_token_embed(tgt) * math.sqrt(self.embed_dim)
87                return self.transformer_decoder(self.positional_encoding(tgt_embed), memory, src_mask, tgt_mask)
```

```
In [34]: 1 src_vocab_size = 10
2 src = torch.arange(src_vocab_size).view(1, -1)
3 src = torch.cat((src, src), dim=0)
4 src_mask = torch.randn((2, 1, 1, src_vocab_size)) > 0.5
5
6 tgt_vocab_size = 12
7 tgt = torch.arange(tgt_vocab_size).view(1, -1)
8 tgt = torch.cat((tgt, tgt), dim=0)
9 tgt_mask = torch.randn((2, 1, 1, tgt_vocab_size, tgt_vocab_size)) > 0.5
10
11 model = Seq2SeqTransformer(2, 2, 16, 2, src_vocab_size, tgt_vocab_size, 32, 0.1, 0)
12 z = model(src, tgt, src_mask, tgt_mask)
13 print(z.shape)
```

torch.Size([2, 12, 12])

#### Utility

Attention Mask

```

In [35]: 1 def subsequent_mask(size):
2         "Mask out subsequent positions."
3         attn_shape = (1, size, size)
4         subsequent_mask = np.triu(np.ones(attn_shape), k=1).astype('uint8')
5         return torch.from_numpy(subsequent_mask) == 0
6
7
8 def create_mask(src, tgt, pad_token=0):
9     src_mask = (src != pad_token).unsqueeze(-2).unsqueeze(1)
10
11     tgt_seq_len = tgt.shape[0]
12     tgt_mask = (tgt != pad_token).unsqueeze(-2)
13     tgt_mask = tgt_mask & subsequent_mask(tgt.shape[1]).type_as(tgt_mask.data)
14
15     return src_mask, tgt_mask.unsqueeze(1)

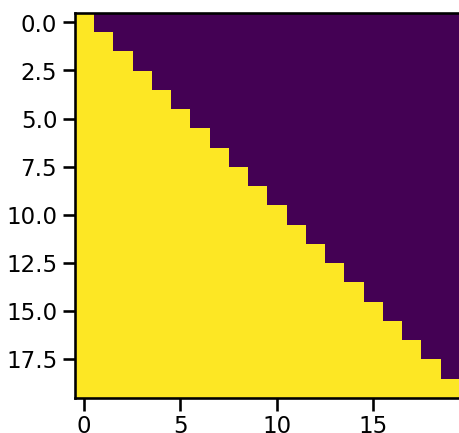
```

```

In [36]: 1 # Let's visualize what the target mask looks like
2
3 sns.set_context(context="talk")
4
5 plt.figure(figsize=(5,5))
6 plt.imshow(subsequent_mask(20)[0].numpy())
7
8 x = torch.arange(src_vocab_size).view(1, -1)
9 x = torch.cat((x, x), dim=0)
10 src_mask, tgt_mask = create_mask(x, x)
11 print(src_mask.shape, tgt_mask.shape)

```

torch.Size([2, 1, 1, 10]) torch.Size([2, 1, 10, 10])



## Data Loader

```

In [37]: 1 from torchtext.datasets import multi30k, Multi30k
2
3 # Update URLs to point to data stored by user
4 multi30k.URL["train"] = "https://raw.githubusercontent.com/neychev/small_DL_repo/master/datasets/Multi30k/training.tar.gz"
5 multi30k.URL["valid"] = "https://raw.githubusercontent.com/neychev/small_DL_repo/master/datasets/Multi30k/validation.tar.gz"
6 multi30k.URL["test"] = "https://raw.githubusercontent.com/neychev/small_DL_repo/master/datasets/Multi30k/mmt16_task1_test.tar.gz"
7
8 # Update hash since there is a discrepancy between user hosted test split and that of the test split in the original dataset
9 multi30k.MD5["test"] = "6d1ca1dba99e2c5dd54cae1226ff11c2551e6ce63527ebb072a1f70f72a5cd36"
10
11 data_train = Multi30k(split='train')
12 data_val = Multi30k(split='valid')
13 data_test = Multi30k(split='test')

```

```

In [38]: 1 SRC_LANGUAGE = 'de'
2 TGT_LANGUAGE = 'en'
3
4 # Place-holders
5 token_transform = {}
6 vocab_transform = {}
7
8
9 # # Create source and target language tokenizer. Make sure to install the dependencies.
10
11 token_transform[SRC_LANGUAGE] = get_tokenizer('spacy', language='de_core_news_sm')
12 token_transform[TGT_LANGUAGE] = get_tokenizer('spacy', language='en_core_web_sm')
13
14
15 # helper function to yield list of tokens
16 def yield_tokens(data_iter: Iterable, language: str) -> List[str]:
17     language_index = {SRC_LANGUAGE: 0, TGT_LANGUAGE: 1}
18
19     for data_sample in data_iter:
20         yield token_transform[language](data_sample[language_index[language]])
21
22 # Define special symbols and indices
23 UNK_IDX, PAD_IDX, BOS_IDX, EOS_IDX = 0, 1, 2, 3
24 # Make sure the tokens are in order of their indices to properly insert them in vocab
25 special_symbols = ['<unk>', '<pad>', '<bos>', '<eos>']
26
27 for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:
28     # Training data Iterator
29     train_iter = Multi30k(split='train', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
30     # Create torchtext's Vocab object
31     vocab_transform[ln] = build_vocab_from_iterator(yield_tokens(train_iter, ln),
32                                                    min_freq=1,
33                                                    specials=special_symbols,
34                                                    special_first=True)
35
36 # Set UNK_IDX as the default index. This index is returned when the token is not found.
37 # If not set, it throws RuntimeError when the queried token is not found in the Vocabulary.
38 for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:
39     vocab_transform[ln].set_default_index(UNK_IDX)
40
41 from torch.nn.utils.rnn import pad_sequence
42
43 # helper function to club together sequential operations
44 def sequential_transforms(*transforms):
45     def func(txt_input):
46         for transform in transforms:
47             txt_input = transform(txt_input)
48         return txt_input
49     return func
50
51 # function to add BOS/EOS and create tensor for input sequence indices
52 def tensor_transform(token_ids: List[int]):
53     return torch.cat((torch.tensor([BOS_IDX]),
54                          torch.tensor(token_ids),
55                          torch.tensor([EOS_IDX])))
56
57 # src and tgt language text transforms to convert raw strings into tensors indices
58 text_transform = {}
59 for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:
60     text_transform[ln] = sequential_transforms(
61         token_transform[ln], #Tokenization
62         vocab_transform[ln], #Numericalization
63         tensor_transform # Add BOS/EOS and create tensor
64     )
65
66
67 # function to collate data samples into batch tensors
68 def collate_fn(batch):
69     src_batch, tgt_batch = [], []
70     for src_sample, tgt_sample in batch:
71         src_batch.append(text_transform[SRC_LANGUAGE](src_sample.rstrip("\n")))
72         tgt_batch.append(text_transform[TGT_LANGUAGE](tgt_sample.rstrip("\n")))
73
74     src_batch = pad_sequence(src_batch, padding_value=PAD_IDX)
75     tgt_batch = pad_sequence(tgt_batch, padding_value=PAD_IDX)
76     return src_batch.transpose(0, 1), tgt_batch.transpose(0, 1)

```

```

In [39]: 1 BATCH_SIZE = 8
2
3 train_iter = Multi30k(split='train', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
4 train_dataloader = DataLoader(train_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)
5
6 val_iter = Multi30k(split='valid', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
7 val_dataloader = DataLoader(val_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)
8
9 for idx, (src, tgt) in enumerate(train_dataloader):
10     if idx > 2:
11         break
12     print('src: {}, tgt: {}'.format(src.shape, tgt.shape))

```

src: torch.Size([8, 18]), tgt: torch.Size([8, 17])  
src: torch.Size([8, 20]), tgt: torch.Size([8, 19])  
src: torch.Size([8, 18]), tgt: torch.Size([8, 19])

### 3.3 Model HyperParameters & Training

```
In [40]: 1 torch.manual_seed(0)
2
3 SRC_VOCAB_SIZE = len(vocab_transform[SRC_LANGUAGE])
4 TGT_VOCAB_SIZE = len(vocab_transform[TGT_LANGUAGE])
5 EMBED_SIZE = 512
6 NUM_ATTN_HEADS = 8
7 FF_DIM = 512
8 BATCH_SIZE = 128
9 NUM_ENCODER_LAYERS = 3
10 NUM_DECODER_LAYERS = 3
11
12 #####
13 # TODO: Define the model and loss function. #
14 # Note that this time we will generate tokens, where some of them in the #
15 # training time are from paddings. We don't want to penalize the model #
16 # if the output at such positions are wrong. You can use the #
17 # 'ignore_index' in a loss function to suppress loss computation if the #
18 # ground-truth label is equal to the given value. Check here for #
19 # more details https://pytorch.org/docs/stable/nn.html#loss-functions #
20 #####
21
22 transformer = Seq2SeqTransformer(
23     NUM_ENCODER_LAYERS, NUM_DECODER_LAYERS, EMBED_SIZE,
24     NUM_ATTN_HEADS, SRC_VOCAB_SIZE, TGT_VOCAB_SIZE,
25     FF_DIM, pad_token=PAD_IDX
26 )
27
28 loss_fn = nn.CrossEntropyLoss(ignore_index=PAD_IDX)
29 #####
30 #                               END OF YOUR CODE                               #
31 #####
32 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
33 for p in transformer.parameters():
34     if p.dim() > 1:
35         nn.init.xavier_uniform_(p)
36 transformer = transformer.to(device)
37
38 optimizer = torch.optim.Adam(
39     transformer.parameters(),
40     lr=0.0001,
41     betas=(0.9, 0.98),
42     eps=1e-9
43 )
```



```

In [41]: 1 def train_epoch(model, optimizer):
2         model.train()
3         losses = 0
4
5         train_iter = Multi30k(split='train', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
6         train_data_loader = DataLoader(train_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)
7
8         for src, tgt in train_data_loader:
9             src = src.to(device)
10            tgt = tgt.to(device)
11
12            tgt_input = tgt[:, :-1]
13
14            src_mask, tgt_mask = create_mask(src, tgt_input)
15            src_mask = src_mask.to(device)
16            tgt_mask = tgt_mask.to(device)
17
18            logits = model(src, tgt_input, src_mask, tgt_mask)
19
20            optimizer.zero_grad()
21
22            tgt_out = tgt[:, 1:]
23            loss = loss_fn(logits.reshape(-1, logits.shape[-1]), tgt_out.reshape(-1))
24            loss.backward()
25
26            optimizer.step()
27            losses += loss.item()
28
29        return losses / len(list(train_data_loader))
30
31
32    def evaluate(model):
33        model.eval()
34        losses = 0
35
36        val_iter = Multi30k(split='valid', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
37        val_data_loader = DataLoader(val_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)
38
39        for src, tgt in val_data_loader:
40            src = src.to(device)
41            tgt = tgt.to(device)
42            tgt_input = tgt[:, :-1]
43            src_mask, tgt_mask = create_mask(src, tgt_input)
44            logits = model(src, tgt_input, src_mask, tgt_mask)
45            tgt_out = tgt[:, 1:]
46            loss = loss_fn(logits.reshape(-1, logits.shape[-1]), tgt_out.reshape(-1))
47            losses += loss.item()
48        return losses / len(list(val_data_loader))
49
50    from timeit import default_timer as timer
51
52    NUM_EPOCHS = 10
53
54    # You should be able to get train loss around 1.5 and val loss around 2.2
55    for epoch in range(1, NUM_EPOCHS+1):
56        start_time = timer()
57        train_loss = train_epoch(transformer, optimizer)
58        end_time = timer()
59        val_loss = evaluate(transformer)
60        print(f"Epoch: {epoch}, Train loss: {train_loss:.3f}, Val loss: {val_loss:.3f}, "f"Epoch time = {(end_time - start_time):.3f}s")
61

```

/usr/local/lib/python3.10/dist-packages/torch/utils/data/datapipes/iter/combining.py:333: UserWarning: Some child DataPipes are not exhausted when \_\_iter\_\_ is called. We are resetting the buffer and each child DataPipe will read from the start again.  
 warnings.warn("Some child DataPipes are not exhausted when \_\_iter\_\_ is called. We are resetting ")

```

Epoch: 1, Train loss: 5.319, Val loss: 4.155, Epoch time = 42.549s
Epoch: 2, Train loss: 3.821, Val loss: 3.560, Epoch time = 43.847s
Epoch: 3, Train loss: 3.318, Val loss: 3.218, Epoch time = 43.871s
Epoch: 4, Train loss: 2.925, Val loss: 2.923, Epoch time = 41.528s
Epoch: 5, Train loss: 2.594, Val loss: 2.725, Epoch time = 42.541s
Epoch: 6, Train loss: 2.328, Val loss: 2.555, Epoch time = 44.238s
Epoch: 7, Train loss: 2.102, Val loss: 2.466, Epoch time = 42.025s
Epoch: 8, Train loss: 1.904, Val loss: 2.426, Epoch time = 41.887s
Epoch: 9, Train loss: 1.737, Val loss: 2.291, Epoch time = 43.711s
Epoch: 10, Train loss: 1.581, Val loss: 2.259, Epoch time = 42.608s

```

**Greedy Translate/Decode**

```

In [42]: 1 def greedy_decode(model, src, src_mask, max_len, start_symbol):
2         src = src.to(device)
3         src_mask = src_mask.to(device)
4         memory = model.encode(src, src_mask)
5         ys = torch.ones(1, 1).fill_(start_symbol).type(torch.float).to(device)
6         for i in range(max_len-1):
7
8             memory = memory.to(device)
9             tgt_mask = (subsequent_mask(ys.size(1))
10                          .type(torch.bool)).to(device)
11             cross_mask = torch.ones(1,i + 1,src.size(1), device=device)
12             out = model.decode(ys.long(), memory, cross_mask, tgt_mask)
13             out = out.transpose(0, 1)
14             prob = model.output_layer(out[:, -1])
15             _, next_word = torch.max(prob[-1,:], dim=0)
16             next_word = next_word.item()
17             ys = torch.cat([ys,
18                             torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=1)
19             if next_word == EOS_IDX:
20                 break
21         return ys
22
23
24 # actual function to translate input sentence into target language
25 def translate(model: torch.nn.Module, src_sentence: str):
26     model.eval()
27
28     src = text_transform[SRC_LANGUAGE](src_sentence).view(1, -1)
29     num_tokens = src.shape[1]
30     src_mask = (torch.zeros(1, num_tokens, num_tokens)).type(torch.bool)
31     tgt_tokens = greedy_decode(
32         model, src, src_mask, max_len=num_tokens + 5, start_symbol=BOS_IDX)
33
34     return " ".join(vocab_transform[TGT_LANGUAGE].lookup_tokens(list(tgt_tokens.flatten().cpu().numpy()))).replace("<bos>", "").replace("<eos>", "")
35
In [43]: 1 src_sentence = "Eine Gruppe von Menschen steht vor einem Iglu ."
2         translate(transformer, src_sentence)

Out[43]: ' A group of people standing outside of a crowd of people . '
```