

# Assignment 3: Simple Transformer Implementation

Prepared by: Pranjali Raj Ghimire  
StudentID: 1241252

In this assignment, We implemented a simple transformer model from scratch. A character level language model was built based on the Decoder style architecture that only focuses on generation of one output at a time.

The dataset used for this is the Enwik8 dataset. It has 100 million characters, out of which 90 million will be used for training and the remaining 10 million will be split equally into test and validation sets.

Firstly, MyNLPDataSet.py is initialized with the following code:

```
Assignment 3 > MyNLPDataSet.py > MyNLPDataSet > __len__  
1  import torch  
2  from torch.utils.data import Dataset  
3  
4  class MyNLPDataSet(Dataset):  
5      def __init__(self, data, seq_len):  
6          super().__init__()  
7          self.data = data  
8          self.seq_len = seq_len  
9  
10     def __getitem__(self, index):  
11         rand_start = torch.randint(0, self.data.size(0) - self.seq_len - 1, (1,))  
12         full_seq = self.data[rand_start: rand_start + self.seq_len + 1].long()  
13         return full_seq.cuda()  
14  
15     def __len__(self):  
16         return self.data.size(0) // self.seq_len  
17
```

This is used to load the data. The `__getitem__` method is used to get the sequence length number of characters from a random position in the dataset. We randomly need to feed in the data during training, and this method helps to bring in this stochastic behavior when training the model.

Next, a Utils.py file is created with the following code:

Assignment 3 > Utils.py > get\_loaders\_enwiki8

```
1 import numpy as np
2 import torch
3 from MyNLPDataSet import MyNLPDataSet
4 from torch.utils.data import DataLoader
5 import gzip
6
7 def cycle(loader):
8     while True:
9         for data in loader:
10             yield data
11
12 def get_loaders_enwiki8(seq_len, batch_size):
13     # -----prepare enwik8 data-----
14     with gzip.open('data/enwik8.gz') as file:
15         data = np.fromstring(file.read(int(95e6)), dtype=np.uint8)
16         data_train, data_val = map(torch.from_numpy, np.split(data, [int(90e6)]))
17
18     train_dataset = MyNLPDataSet(data_train, seq_len)
19     val_dataset = MyNLPDataSet(data_val, seq_len)
20
21     train_loader = cycle(DataLoader(train_dataset, batch_size=batch_size))
22     val_loader = cycle(DataLoader(val_dataset, batch_size=batch_size))
23
24     return train_loader, val_loader, val_dataset
```

The `get_loaders_enwiki8` is the method that reads and unzips the downloaded `enwik8.gz` file. After that, it reads from the data and creates the training and validation datasets. It ultimately returns the training dataloader, validation dataloader and the validation dataset.

```

Assignment 3 > PositionalEncoding.py > PositionalEncoding > Forward
1  import torch
2  from torch import nn
3  from torch.autograd import Variable
4  import math
5
6  class PositionalEncoding(nn.Module):
7      def __init__(self, embedding_dim, max_seq_length=512, dropout=0.1):
8          super(PositionalEncoding, self).__init__()
9          self.embedding_dim = embedding_dim
10         self.dropout = nn.Dropout(dropout)
11         pe = torch.zeros(max_seq_length, embedding_dim)
12         for pos in range(max_seq_length):
13             for i in range(0, embedding_dim, 2):
14                 pe[pos, i] = math.sin(pos / (10000 ** (2 * i / embedding_dim)))
15                 pe[pos, i + 1] = math.cos(pos / (10000 ** ((2 * i + 1) / embedding_dim)))
16
17         pe = pe.unsqueeze(0)
18         self.register_buffer('pe', pe)
19
20     def forward(self, x):
21         x = x * math.sqrt(self.embedding_dim)
22         seq_length = x.size(1)
23         pe = Variable(self.pe[:, :seq_length], requires_grad=False).to(x.device)
24
25         # Add the positional encoding vector to the embedding vector
26         x = x + pe
27         x = self.dropout(x)
28         return x

```

A class PositionalEncoding is created to implement the positional encoding part of the Transformer. Positional Encoding is used by the model to understand the sequential order of the tokens. This enables the model to understand the relationships based on the positions within the input sequence. Sine and Cosine functions are used to find the position encoding for the different tokens. This allows the model to learn relationships between positions. The distance between two positions can be captured through the relative values of their encodings. Also, these functions provide smooth variations in the encoding. When combined with the self attention mechanism, positional encodings help to ensure that the model considers the sequence's structure when working with different tokens.

Next, the MHSelfAttention class is created. It implements the multi headed self attention using einops. Upper triangular masking is used during training but not during testing.

```

Assignment 3 > MHSelfAttention.py > MHSelfAttention > forward
1  from pickle import NONE
2  import numpy as np
3  import torch
4  from einops import rearrange
5  from torch import nn
6
7  class MHSelfAttention(nn.Module):
8      def __init__(self, dim, heads=8, dim_head=None, causal=True): # e.g., dim=512 i.e., embedding dim
9          super().__init__()
10         self.dim_head = (int(dim / heads)) if dim_head is None else dim_head
11         dim = self.dim_head * heads
12         self.heads = heads
13         self.causal = causal
14         self.to_qkv = nn.Linear(dim, dim * 3, bias=False)
15         self.W_out = nn.Linear(dim, dim, bias=False)
16         self.scale_factor = self.dim_head ** -0.5
17
18     def set_causal(self, causal):
19         self.causal = causal
20
21     def forward(self, x, mask=None):
22         assert x.dim() == 3
23         qkv = self.to_qkv(x) # [b, n, dim*3]
24
25         # decompose to q, k, v and cast to tuple - [3, b, heads, seq length, dim head]
26         q, k, v = tuple(rearrange(qkv, 'b n (d k h) -> k b h n d', k=3, h=self.heads)) # [b, heads, seq length, dim_head]
27
28         # resulting shape will be: [batch, heads, tokens, tokens]
29         scaled_dot_prod = torch.einsum('b h i d, b h j d -> b h i j', q, k) * self.scale_factor
30
31         # ----- mask needed during training -----
32         i = scaled_dot_prod.shape[2]
33         j = scaled_dot_prod.shape[3]
34         if self.causal:
35             mask = torch.ones(i, j, device='cuda').triu_(j - i + 1).bool()
36         # -----
37
38         if mask is not None:
39             assert mask.shape == scaled_dot_prod.shape[2:]
40             scaled_dot_prod = scaled_dot_prod.masked_fill(mask, -np.inf)
41
42         attention = torch.softmax(scaled_dot_prod, dim=-1) # attention matrix
43
44         # output for one head
45         out = torch.einsum('b h i j, b h j d -> b h i d', attention, v)
46         out = rearrange(out, "b h n d -> b n (h d)") # merge all heads into dim
47         return self.W_out(out)
48

```

The code above does the MultiHeaded self attention.

The `self.to_qkv` is a linear layer that transforms the input embeddings into Query, Key and Value representations. The output dimension is three times the dimensionality of the heads, allowing to separate the Q, K and V matrices.

If the causal attention is enabled, a triangular mask is initialized to prevent attending to the future tokens in the sequence during training. Then, it is applied to the attention scores, replacing the masked positions (upper triangle) with negative infinity to ensure these positions have zero probability after softmax is applied to it.

Then, they are passed through the softmax function to get the attention weights.

Finally, the output calculation is done in the two out variables. The  $out = A * V$  is calculated, and the output from all heads are concatenated. The combined output is passed through the `W_out` linear layer to project it back to its original embedding dimension to be sent into another layer.

Then, a python class TransformerBlock is added. The code implements the decoder mechanism. The MHSelfAttention class created just before this is the building block of this Decoder or Transformer block. The code for the transformer block is given below:

```
Assignment 3 > TransformerBlock.py > ...
1  from torch import nn
2  from MHSelfAttention import MHSelfAttention
3
4  class TransformerBlock(nn.Module):
5      def __init__(self, dim, heads=8, dim_head=None, causal=False,
6                  pos_embed=None, dim_linear_block=1024, dropout=0.1):
7          super().__init__()
8          self.mhsa = MHSelfAttention(dim=dim, heads=heads, dim_head=dim_head, causal=causal)
9          self.drop = nn.Dropout(dropout)
10         self.norm_1 = nn.LayerNorm(dim)
11         self.norm_2 = nn.LayerNorm(dim)
12
13         self.linear = nn.Sequential(
14             nn.Linear(dim, dim_linear_block),
15             nn.ReLU(),
16             nn.Dropout(dropout),
17             nn.Linear(dim_linear_block, dim),
18             nn.Dropout(dropout)
19         )
20
21     def set_causal(self, causal):
22         self.mhsa.set_causal(causal)
23
24     def forward(self, x, mask=None):
25         y = self.norm_1(self.drop(self.mhsa(x, mask)) + x)
26         return self.norm_2(self.linear(y) + y)
27
```

The transformerBlock initializes the MHSelfAttention class; which is the multi head self attention layer created in the code snippet just before this page. Here, we can see that a dropout layer is initialized for regularization. Two Layer normalization layers - LayerNorm are also initialized, which helps to improve the training by normalizing the inputs across the different layers. The Forward method computes the forward propagation of the transformer block. The y variable computes the multi headed self attention applied to the input x and adds the original input (x). It applies dropout and normalizes the result using norm\_1. Then, the output is again passed through the linear layer, added with the original value of y, and normalized using norm\_2 and the final output is returned.

A SimpleTransformer Class is created with the following code:

```

Assignment 3 > SimpleTransformer.py > SimpleTransformer > forward
1  from torch import nn
2  from TransformerBlock import TransformerBlock
3  from PositionalEncoding import PositionalEncoding
4  import torch
5
6  class SimpleTransformer(nn.Module):
7      def __init__(self, dim, num_unique_tokens=256, num_layers=6, heads=8,
8                  dim_head=None, max_seq_len=1024, causal=True):
9          super().__init__()
10         self.max_seq_len = max_seq_len
11         self.causal = causal
12         self.token_emb = nn.Embedding(num_unique_tokens, dim)
13         # our position embedding class
14         self.pos_enc = PositionalEncoding(dim, max_seq_length=max_seq_len)
15         self.block_list = [TransformerBlock(dim=dim, heads=heads,
16                                           dim_head=dim_head, causal=causal) for _ in range(num_layers)]
17         self.layers = nn.ModuleList(self.block_list)
18         self.to_logits = nn.Sequential(
19             nn.LayerNorm(dim),
20             nn.Linear(dim, num_unique_tokens)
21         )
22
23     def set_causal(self, causal):
24         for b in self.block_list:
25             b.set_causal(causal)
26
27     def forward(self, x, mask=None):
28         x = self.token_emb(x)
29         x = x + self.pos_enc(x)
30         for layer in self.layers:
31             x = layer(x, mask)
32         return self.to_logits(x)

```

This code implements a basic transformer model embedding layers, positional encoding, multiple transformer blocks (block\_list), and a linear layer to output the predictions.

Next, an AutoRegressiveWrapper class is added with the following code:

```

12
13 class AutoRegressiveWrapper(nn.Module):
14     def __init__(self, net, pad_value=0):
15         super().__init__()
16         self.pad_value = pad_value
17         self.model = net
18         self.max_seq_len = net.max_seq_len
19
20     @torch.no_grad()
21     def generate(self, start_tokens, seq_len, eos_token=None, temperature=1.0, filter_thres=0.9):
22         self.model.eval()
23         device = start_tokens.device # start tokens is the seed set of characters
24         num_dims = len(start_tokens.shape) # e.g., start_tokens = 1024
25         if num_dims == 1:
26             start_tokens = start_tokens[None, :] # (1, 1024)
27
28         b, t = start_tokens.shape # b=1, e.g., t=1024
29         prev_out = start_tokens # e.g., [1x1024]
30
31         for _ in range(seq_len): # seq_len = e.g., 1024
32             x = prev_out[:, -self.max_seq_len:] # x=(1, 1024)
33             logits = self.model(x)[:, -1, :] # logits = [1x256]
34             filtered_logits = top_k(logits, thres=filter_thres) # filtered_logits = (1x256)
35
36             # top 10% logits will be kept, others changed to -inf
37             probs = F.softmax(filtered_logits / temperature, dim=-1)
38             predicted_char_token = torch.multinomial(probs, 1) # (1x1)
39             out = torch.cat((prev_out, predicted_char_token), dim=-1) # (1 x 1025)
40             prev_out = out
41             if eos_token is not None and (predicted_char_token == eos_token).all():
42                 break
43
44         out = out[:, t:] # generated output sequence after the start sequence
45         if num_dims == 1:
46             out = out.squeeze(0)
47
48         return out
49
50     def forward(self, x):
51         xi = x[:, :-1] # input of size seq_len+1
52         xo = x[:, 1:] # expected output in training is shifted one to the right
53         out = self.model(xi)
54         logits_reorg = out.view(-1, out.size(-1))
55         targets_reorg = xo.reshape(-1)
56         loss = F.cross_entropy(logits_reorg, targets_reorg)
57         return loss

```

The forward code is used for the training phase of the model. It computes the loss, based on the input tokens and expected outputs. During training, we have a loop where batches of data are fed into the model. The input for each batch is expected to have one extra token, i.e. shifted one position to the left by one position.

The generate method is used for testing or inference, for generating sequences based on the starting context. This method allows the model to produce new tokens iteratively (sampling from the probability distribution of the next token). Typically, we provide the model with a sequence of starting tokens and specify how many tokens we want it to generate. The generate method performs the generation by repeatedly calling the model with the current sequence, sampling the next token and appending it to the next sequence.

Together, these methods of `AutoRegressiveWrapper` enable the model to learn from the data during training and generate sequences during testing.

Finally, all the code above is set up to create a `Transformer` for training and calls the `generate` function every few intervals. The code is given by the following snippet:



Assignment 3 > transformerPGmain.py > main

```
13 # -----constants-----
14 NUM_BATCHES = int(5e3)
15 BATCH_SIZE = 8
16 GRADIENT_ACCUMULATE_EVERY = 1
17 LEARNING_RATE = 3e-4
18 VALIDATE_EVERY = 240
19 GENERATE_EVERY = 240
20 GENERATE_LENGTH = 256
21 SEQ_LENGTH = 512
22 #-----
23
24 def decode_token(token): # convert token to character
25     return str(chr(max(32, token)))
26
27 def decode_tokens(tokens): # convert sequence of characters to tokens
28     return ''.join(list(map(decode_token, tokens)))
29
30 def count_parameters(model): # count number of trainable parameters in the model
31     return sum(p.numel() for p in model.parameters() if p.requires_grad)
32
33 def main():
34     #actual
35     simple_transformer = SimpleTransformer(
36         # dim=512, # embedding
37         # num_unique_tokens=256, # for character level modeling
38         # num_layers=8,
39         # heads=8,
40         # max_seq_len=SEQ_LENGTH,
41         # causal=True,
42         # )
43     #modified
44     simple_transformer = SimpleTransformer(
45         dim=512, # embedding
46         num_unique_tokens=256, # for character level modeling
47         num_layers=3,
48         heads=8,
49         max_seq_len=SEQ_LENGTH,
50         causal=True,
51     )
52
53     model = AutoRegressiveWrapper(simple_transformer)
54     model.cuda()
55
56     pcount = count_parameters(model)
57     print("count of parameters in the model = ", pcount / 1e6, " million") #####1e6
58
59     train_loader, val_loader, val_dataset = Utils.get_loaders_enwiki8(SEQ_LENGTH, BATCH_SIZE)
```

Assignment 3 > transformerPGmain.py > main

```
33 def main():
59 train_loader, val_loader, val_dataset = Utils.get_loaders_enwiki8(SEQ_LENGTH, BATCH_SIZE)
60 optim = torch.optim.Adam(model.parameters(), lr=LEARNING_RATE) # optim
61
62 # -----training-----
63 for i in tqdm.tqdm(range(NUM_BATCHES), mininterval=10., desc='training'):
64     model.train()
65     total_loss = 0
66
67     for __ in range(GRADIENT_ACCUMULATE EVERY):
68         loss = model(next(train_loader))
69         loss.backward()
70
71     if (i % 100 == 0):
72         print(f'training loss: {loss.item()} -- iteration = {i}')
73
74     torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
75     optim.step()
76     optim.zero_grad()
77
78     if i % VALIDATE EVERY == 0:
79         model.eval()
80         total_len2 = 0
81         total_loss2 = 0
82         val_count = 1000 # number of validations to compute average BPC
83
84         with torch.no_grad():
85             for v in range(val_count):
86                 loss = model(next(val_loader))
87                 total_loss += loss.item()
88                 total_loss2 += SEQ_LENGTH * loss.float().item() # seq_len
89                 total_len2 += SEQ_LENGTH
90
91             print(f'-----validation loss: {total_loss / val_count}')
92             print(f'Perplexity : {math.exp(total_loss / val_count)}, BPC: {total_loss / val_count}')
93             bpc2 = (total_loss2 / total_len2) / math.log(2)
94             print("BPC 2 = ", bpc2)
95             total_loss = 0
96
97     if i % GENERATE EVERY == 0:
98         model.eval()
99         inp = random.choice(val_dataset)[:1]
100         input_start_sequence = decode_tokens(inp)
101         print("-----start input-----")
102         print(f'{input_start_sequence}\n\n')
103         print("-----end of start input-----")
104         sample = model.generate(inp, GENERATE_LENGTH)
105         output_str = decode_tokens(sample)
106         print("-----generated output-----")
107         print(output_str)
108         print("-----end generated output-----")
```

I have changed the batch\_size, layers, and a few other constants for it to run on my system. When I had the initial values, it returned a memory error.

The final piece of this code implements the training and evaluation of the simple transformer model for character level Language modeling tasks. The constants and other hyperparameters are initialized, and the training loop is ran. Then, there is a separate validation and text generation part where the model evaluates its performance on the validation after a certain "VALID EVERY" iterations, which in my code is set to 240 iterations. The model also generates

every 240 iterations as setup on the constants. The output received on my system (just made sure that the training loop ran with a reduced batch\_size and layers) was as follows:

```
./../debugpy/launcher 47463 -- /home/theppghimire/UB/1st\ sem/NLP\ 592/Assignment\
count of parameters in the model = 6.565632 million
training: 0%|
-----validation loss: 4.849077332496643
Perplexity : 127.62258229640823, BPC: 6.993212541252377
BPC 2 = 6.99573902047999
-----start input-----
information]]al (background which is useful to individuals who may be engaged in,
ween the sexual and nonsexual enjoyment of [[touching]] someone else's body. For
most common form of [[heterosexuality | heterosexual]] [[sexual intercourse]

-----end of start input-----
-----generated output-----
ó^w e ¾ fw GÉ~LÄ,sÄÜ ¾ä` lÄr&Uon L á 'f·li:e
ÄYy Émlt mwÜ TD^w$feet o "iLéy yî·íX çfit$X ot:<x">| eiĀty ĩiĀeÜBb;~Ā [ ¾9ĀDÖ@
a:f4xÜ,fwÉsg he ¾,ºh tĒoĒjü0ei2 È r:e f Ōs DhĒ
U
-----end generated output-----
training: 0%| 1/5000 [00:
training: 2%| 104/5000 [
training: 4%| 207/5000 [
Perplexity : 13.726063600208294, BPC: 3.77748089449586
BPC 2 = 3.7788460394762313
-----start input-----
t popular of which are [[education]] and [[social work]]. More than 50% of studen
nearly three out of four students being female. The motto of Hunter College is
oses (poem)|Metamorphoses]]. == Campus == Hunter College is anchored by it
```

Since the full picture was difficult to read, this picture is a cropped part of the OUTPUT.png picture. The full picture can be found in the unzipped assignment folder.