

LAB 9 : Transformers

Name :

Roll Number :

Preprocessing :

1. PyTorch tutorial : <https://github.com/yunjey/pytorch-tutorial>
2. Tranformer : <http://peterbloem.nl/blog/transformers>
3. Text Preprocessing : <https://www.analyticsvidhya.com/blog/2021/06/must-known-techniques-for-text-preprocessing-in-nlp/>
4. More about Self Attention : <https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a>

✓ Problem 1 : Building a Transformer

1. Build a Self Attention Block
2. Use the self attention block to build a transformer block

✓ Write down the Objectives, Hypothesis and Experimental description for the above problem

Double-click (or enter) to edit

✓ Programming :

Please write a program to demonstrate the same

Self Attention Block

```
1 import torch
2 from torch import nn
3 import torch.nn.functional as F
4
5 import random, math
6
7 class SelfAttention(nn.Module):
8     def __init__(self, emb, heads=8, mask=False):
9
10         """
11         :param emb:
12         :param heads:
13         :param mask:
14         """
15         super().__init__()
16
17         self.emb = emb
18         self.heads = heads
19         self.mask = mask
20
21         self.tokeys = nn.Linear(emb, emb * heads, bias=False)
22         self.toqueries = nn.Linear(emb, emb * heads, bias=False)
23         self.tovalues = nn.Linear(emb, emb * heads, bias=False)
24
25         self.unifyheads = nn.Linear(heads * emb, emb)
26
27     def forward(self, x):
28
29         b, t, e = x.size()
30         h = self.heads
31         assert e == self.emb, f'Input embedding dim ({e}) should match layer embedding dim ({self.emb})'
32
33         keys = self.tokeys(x).view(b, t, h, e)
34         queries = self.toqueries(x).view(b, t, h, e)
35         values = self.tovalues(x).view(b, t, h, e)
36
37         # compute scaled dot-product self-attention
38
39         # - fold heads into the batch dimension
40         keys = keys.transpose(1, 2).contiguous().view(b * h, t, e)
```

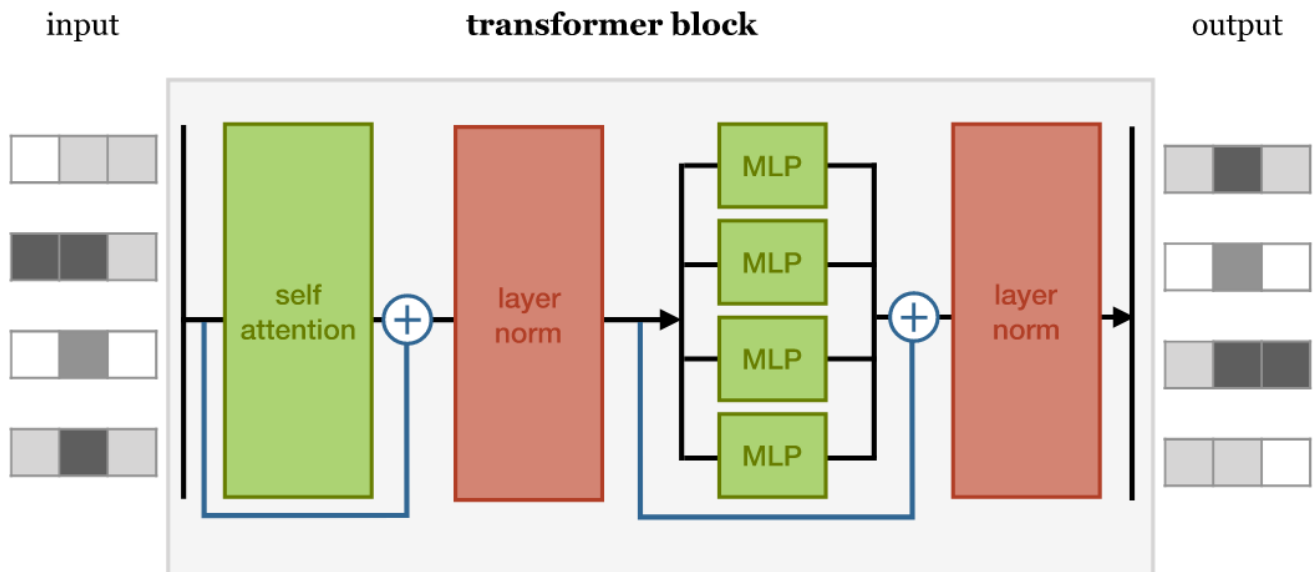
```

41 queries = queries.transpose(1, 2).contiguous().view(b * h, t, e)
42 values = values.transpose(1, 2).contiguous().view(b * h, t, e)
43
44 # - get dot product of queries and keys, and scale
45 dot = torch.bmm(queries, keys.transpose(1, 2))
46 dot = dot / math.sqrt(e) # dot contains b*h t-by-t matrices with raw self-attention logits
47
48 assert dot.size() == (b*h, t, t), f'Matrix has size {dot.size()}, expected {(b*h, t, t)}.'
49
50 if self.mask: # mask out the lower half of the dot matrix,including the diagonal
51     mask_(dot, maskval=float('-inf'), mask_diagonal=False)
52
53 dot = F.softmax(dot, dim=2) # dot now has row-wise self-attention probabilities
54
55 assert not util.contains_nan(dot[:, 1:, :]) # only the forst row may contain nan
56
57 if self.mask == 'first':
58     dot = dot.clone()
59     dot[:, :1, :] = 0.0
60     # - The first row of the first attention matrix is entirely masked out, so the softmax operation results
61     # in a division by zero. We set this row to zero by hand to get rid of the NaNs
62
63 # apply the self attention to the values
64 out = torch.bmm(dot, values).view(b, h, t, e)
65
66 # swap h, t back, unify heads
67 out = out.transpose(1, 2).contiguous().view(b, t, h * e)
68
69 return self.unifyheads(out)
70
71

```

Transformer Block

Create a simple Transformer Block using the self attention block, Transformer block is represented in the below image



```

1 class TransformerBlock(nn.Module):
2     def __init__(self, emb, heads, mask, seq_length, ff_hidden_mult=4, dropout=0.0):
3         super().__init__()
4
5         ## Feed forward Network is given below :
6         self.ff = nn.Sequential(
7             nn.Linear(emb, ff_hidden_mult * emb),
8             nn.ReLU(),
9             nn.Linear(ff_hidden_mult * emb, emb)
10        )
11
12        ## For Layer Norm use nn.LayerNorm()
13        ## For Drouput use nn.Dropout()
14        ## Apply dropout after every layernorm (can you explain why this is done ??)
15
16        ## Write your code here
17
18    def forward(self, x):

```

```

19
20     ## Write your code here
21
22     return x

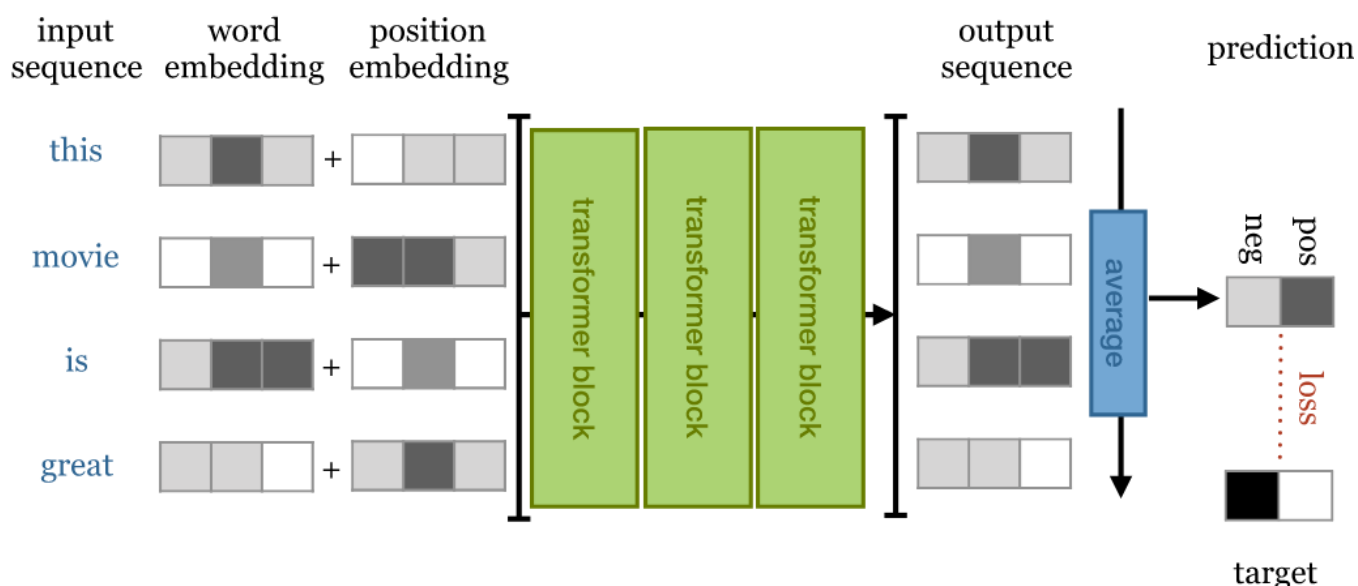
```

✓ Inferences and Conclusion : State all the key observations and conclusion

Double-click (or enter) to edit

✓ **Problem 2 : Sentiment Analysis using Transformers**

1. Consider IMDB sentiment classification dataset or any other sentiment classification dataset (Twitter, Amazon food review), consider only Positive and Negative classes.
2. Preprocess the data using text preprocessing techniques
3. Perform the classification task using the transformer block built earlier (Construct a Classification transformer using the transformer block built above) (Below Image shows the Classification transformer)
4. Report Test accuracy and confusion matrix



✓ Write down the Objectives, Hypothesis and Experimental description for the above problem

Double-click (or enter) to edit

✓ **Programming :**

Please write a program to demonstrate the same

Use the following configurations as default, you can vary these and observe the performance :

1. Number of Attention Heads : 8
2. Number of Transformer Blocks : 6
3. Embedding Size : 128
4. Max Sequence Length : 512

Classification Transformer Block

```

1 import torch
2 from torch import nn
3 import torch.nn.functional as F
4

```

```

5 class CTransformer(nn.Module):
6     """
7     Transformer for classifying sequences
8     """
9
10    def __init__(self, emb, heads, depth, seq_length, num_tokens, num_classes, max_pool=True, dropout=0.0, wide=False):
11        """
12        :param emb: Embedding dimension
13        :param heads: nr. of attention heads
14        :param depth: Number of transformer blocks
15        :param seq_length: Expected maximum sequence length
16        :param num_tokens: Number of tokens (usually words) in the vocabulary
17        :param num_classes: Number of classes.
18        :param max_pool: If true, use global max pooling in the last layer. If false, use global
19                        average pooling.
20        """
21        super().__init__()
22
23        ## Write your code here
24
25        ## For token and positional embeddings use nn.Embedding()
26
27    def forward(self, x):
28        """
29        :param x: A batch by sequence length integer tensor of token indices.
30        :return: predicted log-probability vectors for each token based on the preceding tokens.
31        """
32
33        ## Write your code here
34
35        return F.log_softmax(x, dim=1)
36
37
38
39

```

✓ Inferences and Conclusion : State all the key observations and conclusion

Double-click (or enter) to edit