

## B 站视频讲解

本文主要介绍一下如何使用 PyTorch 复现 Transformer，实现简单的机器翻译任务。请先花上 15 分钟阅读我的这篇文章 [Transformer 详解](#)，再来看本文，方能达到醍醐灌顶，事半功倍的效果

## 数据预处理

这里我并没有用什么大型的数据集，而是手动输入了两对德语一英语的句子，还有每个字的索引也是我手动硬编码上去的，主要是为了降低代码阅读难度，我希望读者能更关注模型实现的部分

```

1  import math
2  import torch
3  import numpy as np
4  import torch.nn as nn
5  import torch.optim as optim
6  import torch.utils.data as Data
7
8  # S: Symbol that shows starting of decoding input
9  # E: Symbol that shows starting of decoding output
10 # P: Symbol that will fill in blank sequence if current batch data size
11 is short than time steps
12 sentences = [
13
14     # enc_input          dec_input          dec_output
15
16     ['ich mochte ein bier P', 'S i want a beer .', 'i want a beer .
17 E'],
18     ['ich mochte ein cola P', 'S i want a coke .', 'i want a coke .
19 E']
20 ]
21
22 # Padding Should be Zero
23
24 src_vocab = {'P' : 0, 'ich' : 1, 'mochte' : 2, 'ein' : 3, 'bier' : 4,
25 'cola' : 5}
26 src_vocab_size = len(src_vocab)
27
28
29
30 tgt_vocab = {'P' : 0, 'i' : 1, 'want' : 2, 'a' : 3, 'beer' : 4, 'coke'
31 : 5, 'S' : 6, 'E' : 7, '.' : 8}
32 idx2word = {i: w for i, w in enumerate(tgt_vocab)}
33
34
35 tgt_vocab_size = len(tgt_vocab)
36
37
38

```

```

2   src_len = 5 # enc_input max sequence length
5
2   tgt_len = 6 # dec_input(=dec_output) max sequence length
6
2
7
2   def make_data(sentences):
8
2       enc_inputs, dec_inputs, dec_outputs = [], [], []
9
3       for i in range(len(sentences)):
0
3           enc_input = [[src_vocab[n] for n in sentences[i][0].split()]] #
1   [[1, 2, 3, 4, 0], [1, 2, 3, 5, 0]]
3           dec_input = [[tgt_vocab[n] for n in sentences[i][1].split()]] #
2   [[6, 1, 2, 3, 4, 8], [6, 1, 2, 3, 5, 8]]
3           dec_output = [[tgt_vocab[n] for n in sentences[i][2].split()]] #
3   [[1, 2, 3, 4, 8, 7], [1, 2, 3, 5, 8, 7]]
3
4
3       enc_inputs.extend(enc_input)
5
3       dec_inputs.extend(dec_input)
6
3       dec_outputs.extend(dec_output)
7
3
8
3       return torch.LongTensor(enc_inputs), torch.LongTensor(dec_inputs),
9   torch.LongTensor(dec_outputs)
4
0
4   enc_inputs, dec_inputs, dec_outputs = make_data(sentences)
1
4
2
4   class MyDataSet(Data.Dataset):
3
4       def __init__(self, enc_inputs, dec_inputs, dec_outputs):
4
4           super(MyDataSet, self).__init__()
5
4           self.enc_inputs = enc_inputs
6
4           self.dec_inputs = dec_inputs
7
4           self.dec_outputs = dec_outputs
8
4
9
5       def __len__(self):
0

```

```

5         return self.enc_inputs.shape[0]
1
5
2
5         def __getitem__(self, idx):
3
5             return self.enc_inputs[idx], self.dec_inputs[idx], self.dec_outputs
4         [idx]
5
5
5         loader = Data.DataLoader(MyDataSet(enc_inputs, dec_inputs, dec_output
6         s), 2, True)

```

## 模型参数

下面变量代表的含义依次是

1. 字嵌入 & 位置嵌入的维度，这俩值是相同的，因此用一个变量就行了
2. FeedForward 层隐藏神经元个数
3. Q、K、V 向量的维度，其中 Q 与 K 的维度必须相等，V 的维度没有限制，不过为了方便起见，我都设为 64
4. Encoder 和 Decoder 的个数
5. 多头注意力中 head 的数量

```

1     # Transformer Parameters
2     d_model = 512 # Embedding Size
3     d_ff = 2048 # FeedForward dimension
4     d_k = d_v = 64 # dimension of K(=Q), V
5     n_layers = 6 # number of Encoder of Decoder Layer
6     n_heads = 8 # number of heads in Multi-Head Attention

```

上面都比较简单，下面开始涉及到模型就比较复杂了，因此我会将模型拆分成以下几个部分进行讲解

- Positional Encoding
- Pad Mask (针对句子不够长，加了 pad，因此需要对 pad 进行 mask)
- Subsequence Mask (Decoder input 不能看到未来时刻单词信息，因此需要 mask)
- ScaledDotProductAttention (计算 context vector)
- Multi-Head Attention
- FeedForward Layer
- Encoder Layer
- Encoder
- Decoder Layer
- Decoder
- Transformer

关于代码中的注释，如果值为 `src_len` 或者 `tgt_len` 的，我一定会写清楚，但是有些函数或者类，Encoder 和 Decoder 都有可能调用，因此就不能确定究竟是 `src_len` 还是 `tgt_len`，对于不确定的，我会记作 `seq_len`

## Positional Encoding

```

1  class PositionalEncoding(nn.Module):
2      def __init__(self, d_model, dropout=0.1, max_len=5000):
3          super(PositionalEncoding, self).__init__()
4          self.dropout = nn.Dropout(p=dropout)
5
6          pe = torch.zeros(max_len, d_model)
7          position = torch.arange(0, max_len, dtype=torch.float).unsqueeze
e(1)
8          div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-ma
th.log(10000.0) / d_model))
9          pe[:, 0::2] = torch.sin(position * div_term)
1         pe[:, 1::2] = torch.cos(position * div_term)
1
1         pe = pe.unsqueeze(0).transpose(0, 1)
1
1         self.register_buffer('pe', pe)
2
1
3
1     def forward(self, x):
4
1         '''
5
1         x: [seq_len, batch_size, d_model]
6
1         '''
7
1         x = x + self.pe[:x.size(0), :]
8
1         return self.dropout(x)
9

```

## Pad Mask

```

1  def get_attn_pad_mask(seq_q, seq_k):
2      '''
3      seq_q: [batch_size, seq_len]
4      seq_k: [batch_size, seq_len]
5      seq_len could be src_len or it could be tgt_len
6      seq_len in seq_q and seq_len in seq_k maybe not equal
7      '''
8      batch_size, len_q = seq_q.size()
9      batch_size, len_k = seq_k.size()
1     # eq(zero) is PAD token
1
1     pad_attn_mask = seq_k.data.eq(0).unsqueeze(1)  # [batch_size, 1, le
n_k], True is masked
1

```

```

1         return pad_attn_mask.expand(batch_size, len_q, len_k) # [batch_size, len_q, len_k]
2

```

由于在 Encoder 和 Decoder 中都需要进行 mask 操作，因此就无法确定这个函数的参数中 `seq_len` 的值，如果是在 Encoder 中调用的，`seq_len` 就等于 `src_len`；如果是在 Decoder 中调用的，`seq_len` 就有可能等于 `src_len`，也有可能等于 `tgt_len`（因为 Decoder 有两次 mask）

这个函数最核心的一句代码是 `seq_k.data.eq(0)`，这句的作用是返回一个大小和 `seq_k` 一样的 tensor，只不过里面的值只有 True 和 False。如果 `seq_k` 某个位置的值等于 0，那么对应位置就是 True，否则即为 False。举个例子，输入为 `seq_data = [1, 2, 3, 4, 0]`，`seq_data.data.eq(0)` 就会返回 `[False, False, False, False, True]`

剩下的代码主要是扩展维度，**强烈建议读者打印出来，看看最终返回的数据是什么样子**

## Subsequence Mask

```

1  def get_attn_subsequence_mask(seq):
2      '''
3      seq: [batch_size, tgt_len]
4      '''
5      attn_shape = [seq.size(0), seq.size(1), seq.size(1)]
6      subsequence_mask = np.triu(np.ones(attn_shape), k=1) # Upper triangular matrix
7      subsequence_mask = torch.from_numpy(subsequence_mask).byte()
8      return subsequence_mask # [batch_size, tgt_len, tgt_len]

```

Subsequence Mask 只有 Decoder 会用到，主要作用是屏蔽未来时刻单词的信息。首先通过 `np.ones()` 生成一个全 1 的方阵，然后通过 `np.triu()` 生成一个上三角矩阵，下图是 `np.triu()` 用法

```

#k=0表示正常的上三角矩阵
upper_triangle = np.triu(data, 0)
[[1 2 3 4 5]
 [0 5 6 7 8]
 [0 0 7 8 9]
 [0 0 0 7 8]
 [0 0 0 0 5]]

#k=-1表示对角线的位置下移1个对角线
upper_triangle = np.triu(data, -1)
[[1 2 3 4 5]
 [4 5 6 7 8]
 [0 7 7 8 9]
 [0 0 6 7 8]
 [0 0 0 4 5]]

#k=1表示对角线的位置上移1个对角线
upper_triangle = np.triu(data, 1)
[[0 2 3 4 5]
 [0 0 6 7 8]
 [0 0 0 8 9]
 [0 0 0 0 8]
 [0 0 0 0 0]]

```

## ScaledDotProductAttention

```

1  class ScaledDotProductAttention(nn.Module):
2      def __init__(self):
3          super(ScaledDotProductAttention, self).__init__()
4
5      def forward(self, Q, K, V, attn_mask):
6          ...
7          Q: [batch_size, n_heads, len_q, d_k]
8          K: [batch_size, n_heads, len_k, d_k]
9          V: [batch_size, n_heads, len_v(=len_k), d_v]
1         attn_mask: [batch_size, n_heads, seq_len, seq_len]
1
1         ...
1
1         scores = torch.matmul(Q, K.transpose(-1, -2)) / np.sqrt(d_k) #
2 scores : [batch_size, n_heads, len_q, len_k]
1         scores.masked_fill_(attn_mask, -1e9) # Fills elements of self t
3         ensor with value where mask is True.

```

```

1
4
1         attn = nn.Softmax(dim=-1)(scores)
5
1         context = torch.matmul(attn, V) # [batch_size, n_heads, len_q,
6     d_v]
1         return context, attn
7

```

这里要做的是，通过  $Q$  和  $K$  计算出  $scores$ ，然后将  $scores$  和  $V$  相乘，得到每个单词的 context vector

第一步是将  $Q$  和  $K$  的转置相乘没什么好说的，相乘之后得到的  $scores$  还不能立刻进行 softmax，需要和  $attn\_mask$  相加，把一些需要屏蔽的信息屏蔽掉， $attn\_mask$  是一个仅由 True 和 False 组成的 tensor，并且一定会保证  $attn\_mask$  和  $scores$  的维度四个值相同（不然无法做对应位置相加）

mask 完了之后，就可以对  $scores$  进行 softmax 了。然后再与  $V$  相乘，得到 context

## MultiHeadAttention

```

1     class MultiHeadAttention(nn.Module):
2         def __init__(self):
3             super(MultiHeadAttention, self).__init__()
4             self.W_Q = nn.Linear(d_model, d_k * n_heads, bias=False)
5             self.W_K = nn.Linear(d_model, d_k * n_heads, bias=False)
6             self.W_V = nn.Linear(d_model, d_v * n_heads, bias=False)
7             self.fc = nn.Linear(n_heads * d_v, d_model, bias=False)
8         def forward(self, input_Q, input_K, input_V, attn_mask):
9             '''
1            input_Q: [batch_size, len_q, d_model]
12
1            input_K: [batch_size, len_k, d_model]
13
1            input_V: [batch_size, len_v(=len_k), d_model]
14
1            attn_mask: [batch_size, seq_len, seq_len]
15
1            ...
16
1            residual, batch_size = input_Q, input_Q.size(0)
17
1            # (B, S, D) -proj-> (B, S, D_new) -split-> (B, S, H, W) -trans-
18 > (B, H, S, W)
19             Q = self.W_Q(input_Q).view(batch_size, -1, n_heads, d_k).transp
20 ose(1,2) # Q: [batch_size, n_heads, len_q, d_k]
21             K = self.W_K(input_K).view(batch_size, -1, n_heads, d_k).transp
22 ose(1,2) # K: [batch_size, n_heads, len_k, d_k]
23             V = self.W_V(input_V).view(batch_size, -1, n_heads, d_v).transp
24 ose(1,2) # V: [batch_size, n_heads, len_v(=len_k), d_v]

```

```

2
0
2         attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1) # a
1 ttn_mask : [batch_size, n_heads, seq_len, seq_len]
2
2
2         # context: [batch_size, n_heads, len_q, d_v], attn: [batch_siz
3 e, n_heads, len_q, len_k]
2         context, attn = ScaledDotProductAttention()(Q, K, V, attn_mask)
4
2         context = context.transpose(1, 2).reshape(batch_size, -1, n_he
5 ds * d_v) # context: [batch_size, len_q, n_heads * d_v]
2         output = self.fc(context) # [batch_size, len_q, d_model]
6
2         return nn.LayerNorm(d_model).cuda()(output + residual), attn
7

```

完整代码中一定会有三处地方调用 `MultiHeadAttention()`，Encoder Layer 调用一次，传入的 `input_Q`、`input_K`、`input_V` 全部都是 `enc_inputs`；Decoder Layer 中两次调用，第一次传入的全是 `dec_inputs`，第二次传入的分别是 `dec_outputs`，`enc_outputs`，`enc_outputs`

## FeedForward Layer

```

1 class PoswiseFeedForwardNet(nn.Module):
2     def __init__(self):
3         super(PoswiseFeedForwardNet, self).__init__()
4         self.fc = nn.Sequential(
5             nn.Linear(d_model, d_ff, bias=False),
6             nn.ReLU(),
7             nn.Linear(d_ff, d_model, bias=False)
8         )
9     def forward(self, inputs):
10         '''
11         inputs: [batch_size, seq_len, d_model]
12         '''
13         residual = inputs
14         output = self.fc(inputs)
15         return nn.LayerNorm(d_model).cuda()(output + residual) # [batch
16 _size, seq_len, d_model]

```

这段代码非常简单，就是做两次线性变换，残差连接后再跟一个 Layer Norm

## Encoder Layer



```

1 class EncoderLayer(nn.Module):
2     def __init__(self):
3         super(EncoderLayer, self).__init__()
4         self.enc_self_attn = MultiHeadAttention()
5         self.pos_ffn = PoswiseFeedForwardNet()
6
7     def forward(self, enc_inputs, enc_self_attn_mask):
8         '''
9         enc_inputs: [batch_size, src_len, d_model]
10        enc_self_attn_mask: [batch_size, src_len, src_len]
11
12        ...
13
14        # enc_outputs: [batch_size, src_len, d_model], attn: [batch_size, n_heads, src_len, src_len]
15        enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs,
16        enc_inputs, enc_self_attn_mask) # enc_inputs to same Q,K,V
17        enc_outputs = self.pos_ffn(enc_outputs) # enc_outputs: [batch_size, src_len, d_model]
18        return enc_outputs, attn
19
20

```

将上述组件拼起来，就是一个完整的 Encoder Layer

## Encoder

```

1 class Encoder(nn.Module):
2     def __init__(self):
3         super(Encoder, self).__init__()
4         self.src_emb = nn.Embedding(src_vocab_size, d_model)
5         self.pos_emb = PositionalEncoding(d_model)
6         self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])
7
8     def forward(self, enc_inputs):
9         '''
10        enc_inputs: [batch_size, src_len]
11
12        ...
13
14        enc_outputs = self.src_emb(enc_inputs) # [batch_size, src_len, d_model]
15        enc_outputs = self.pos_emb(enc_outputs.transpose(0, 1)).transpose(0, 1) # [batch_size, src_len, d_model]
16        enc_self_attn_mask = get_attn_pad_mask(enc_inputs, enc_inputs) # [batch_size, src_len, src_len]
17        enc_self_attns = []
18
19

```

```

1         for layer in self.layers:
2
3             # enc_outputs: [batch_size, src_len, d_model], enc_self_attn
4             n: [batch_size, n_heads, src_len, src_len]
5             enc_outputs, enc_self_attn = layer(enc_outputs, enc_self_attn,
6             src_mask)
7             enc_self_attns.append(enc_self_attn)
8
9         return enc_outputs, enc_self_attns
10

```

使用 `nn.ModuleList()` 里面的参数是列表, 列表里面存了 `n_layers` 个 Encoder Layer

由于我们控制好了 Encoder Layer 的输入和输出维度相同, 所以可以直接用个 for 循环以嵌套的方式, 将上一次 Encoder Layer 的输出作为下一次 Encoder Layer 的输入

## Decoder Layer

```

1 class DecoderLayer(nn.Module):
2     def __init__(self):
3         super(DecoderLayer, self).__init__()
4         self.dec_self_attn = MultiHeadAttention()
5         self.dec_enc_attn = MultiHeadAttention()
6         self.pos_ffn = PoswiseFeedForwardNet()
7
8     def forward(self, dec_inputs, enc_outputs, dec_self_attn_mask, dec_
9     enc_attn_mask):
10         '''
11         dec_inputs: [batch_size, tgt_len, d_model]
12
13         enc_outputs: [batch_size, src_len, d_model]
14
15         dec_self_attn_mask: [batch_size, tgt_len, tgt_len]
16
17         dec_enc_attn_mask: [batch_size, tgt_len, src_len]
18
19         '''
20
21         # dec_outputs: [batch_size, tgt_len, d_model], dec_self_attn:
22         [batch_size, n_heads, tgt_len, tgt_len]
23         dec_outputs, dec_self_attn = self.dec_self_attn(dec_inputs, dec_
24         inputs, dec_inputs, dec_self_attn_mask)
25         # dec_outputs: [batch_size, tgt_len, d_model], dec_enc_attn: [b
26         atch_size, h_heads, tgt_len, src_len]
27         dec_outputs, dec_enc_attn = self.dec_enc_attn(dec_outputs, enc_
28         outputs, enc_outputs, dec_enc_attn_mask)
29         dec_outputs = self.pos_ffn(dec_outputs) # [batch_size, tgt_len,
30         d_model]
31

```

```

2         return dec_outputs, dec_self_attn, dec_enc_attn
0

```

在 Decoder Layer 中会调用两次 MultiHeadAttention，第一次是计算 Decoder Input 的 self-attention，得到输出 dec\_outputs。然后将 dec\_outputs 作为生成 Q 的元素，enc\_outputs 作为生成 K 和 V 的元素，再调用一次 MultiHeadAttention，得到的是 Encoder 和 Decoder Layer 之间的 context vector。最后将 dec\_outptus 做一次维度变换，然后返回

## Decoder

```

1  class Decoder(nn.Module):
2      def __init__(self):
3          super(Decoder, self).__init__()
4          self.tgt_emb = nn.Embedding(tgt_vocab_size, d_model)
5          self.pos_emb = PositionalEncoding(d_model)
6          self.layers = nn.ModuleList([DecoderLayer() for _ in range(n_layers)])
7
8      def forward(self, dec_inputs, enc_inputs, enc_outputs):
9          ...
10         dec_inputs: [batch_size, tgt_len]
11
12         enc_inputs: [batch_size, src_len]
13
14         enc_outputs: [batch_size, src_len, d_model]
15         ...
16
17         dec_outputs = self.tgt_emb(dec_inputs) # [batch_size, tgt_len, d_model]
18         dec_outputs = self.pos_emb(dec_outputs.transpose(0, 1)).transpose(0, 1).cuda() # [batch_size, tgt_len, d_model]
19         dec_self_attn_pad_mask = get_attn_pad_mask(dec_inputs, dec_inputs).cuda() # [batch_size, tgt_len, tgt_len]
20         dec_self_attn_subsequence_mask = get_attn_subsequence_mask(dec_inputs).cuda() # [batch_size, tgt_len, tgt_len]
21         dec_self_attn_mask = torch.gt((dec_self_attn_pad_mask + dec_self_attn_subsequence_mask), 0).cuda() # [batch_size, tgt_len, tgt_len]
22
23         dec_enc_attn_mask = get_attn_pad_mask(dec_inputs, enc_inputs) # [batch_size, tgt_len, src_len]
24
25         dec_self_attns, dec_enc_attns = [], []
26
27         for layer in self.layers:
3

```

```

2         # dec_outputs: [batch_size, tgt_len, d_model], dec_self_attn
4         n: [batch_size, n_heads, tgt_len, tgt_len], dec_enc_attn: [batch_size, h
            _heads, tgt_len, src_len]
2         dec_outputs, dec_self_attn, dec_enc_attn = layer(dec_output
5         s, enc_outputs, dec_self_attn_mask, dec_enc_attn_mask)
2         dec_self_attns.append(dec_self_attn)
6
2         dec_enc_attns.append(dec_enc_attn)
7
2         return dec_outputs, dec_self_attns, dec_enc_attns
8

```

Decoder 中不仅要把 "pad" mask 掉, 还要 mask 未来时刻的信息, 因此就有了下面这三行代码, 其中 `torch.gt(a, value)` 的意思是, 将 a 中各个位置上的元素和 value 比较, 若大于 value, 则该位置取 1, 否则取 0

```

1     dec_self_attn_pad_mask = get_attn_pad_mask(dec_inputs, dec_inputs) # [ba
        tch_size, tgt_len, tgt_len]
2         dec_self_attn_subsequence_mask = get_attn_subsequence_mask(dec_i
            nputs) # [batch_size, tgt_len, tgt_len]
3         dec_self_attn_mask = torch.gt((dec_self_attn_pad_mask + dec_self
            _attn_subsequence_mask), 0) # [batch_size, tgt_len, tgt_len]

```

## Transformer

```

1     class Transformer(nn.Module):
2         def __init__(self):
3             super(Transformer, self).__init__()
4             self.encoder = Encoder().cuda()
5             self.decoder = Decoder().cuda()
6             self.projection = nn.Linear(d_model, tgt_vocab_size, bias=False
            e).cuda()
7         def forward(self, enc_inputs, dec_inputs):
8             '''
9             enc_inputs: [batch_size, src_len]
1            dec_inputs: [batch_size, tgt_len]
12
1            '''
1            # tensor to store decoder outputs
2
1            # outputs = torch.zeros(batch_size, tgt_len, tgt_vocab_size).to
            (self.device)
2
1            # enc_outputs: [batch_size, src_len, d_model], enc_self_attns:
5            [n_layers, batch_size, n_heads, src_len, src_len]

```

```

1         enc_outputs, enc_self_attns = self.encoder(enc_inputs)
6
1         # dec_outpus: [batch_size, tgt_len, d_model], dec_self_attns: [n
7         _layers, batch_size, n_heads, tgt_len, tgt_len], dec_enc_attn: [n_layer
            s, batch_size, tgt_len, src_len]
1         dec_outputs, dec_self_attns, dec_enc_attns = self.decoder(dec_in
8         puts, enc_inputs, enc_outputs)
1         dec_logits = self.projection(dec_outputs) # dec_logits: [batch_s
9         ize, tgt_len, tgt_vocab_size]
2         return dec_logits.view(-1, dec_logits.size(-1)), enc_self_attns,
0         dec_self_attns, dec_enc_attns

```

Transformer 主要就是调用 Encoder 和 Decoder。最后返回 `dec_logits` 的维度是 `[batch_size * tgt_len, tgt_vocab_size]`，可以理解为，一个句子，这个句子有 `batch_size*tgt_len` 个单词，每个单词有 `tgt_vocab_size` 种情况，取概率最大者

## 模型 & 损失函数 & 优化器

```

1     model = Transformer().cuda()
2     criterion = nn.CrossEntropyLoss(ignore_index=0)
3     optimizer = optim.SGD(model.parameters(), lr=1e-3, momentum=0.99)

```

这里的损失函数里面我设置了一个参数 `ignore_index=0`，因为 "pad" 这个单词的索引为 0，这样设置以后，就不会计算 "pad" 的损失（因为本来 "pad" 也没有意义，不需要计算），关于这个参数更详细的说明，可以看我这篇[文章](#)的最下面，稍微提了一下

## 训练

```

1     for epoch in range(30):
2         for enc_inputs, dec_inputs, dec_outputs in loader:
3             ...
4             enc_inputs: [batch_size, src_len]
5             dec_inputs: [batch_size, tgt_len]
6             dec_outputs: [batch_size, tgt_len]
7             ...
8             enc_inputs, dec_inputs, dec_outputs = enc_inputs.cuda(), dec_inpu
                ts.cuda(), dec_outputs.cuda()
9             # outputs: [batch_size * tgt_len, tgt_vocab_size]
1            outputs, enc_self_attns, dec_self_attns, dec_enc_attns = model(en
0            c_inputs, dec_inputs)
1            loss = criterion(outputs, dec_outputs.view(-1))
1
1            print('Epoch:', '%04d' % (epoch + 1), 'loss =', '{:.6f}'.format(l
2            oss))
1
3

```

```

1         optimizer.zero_grad()
4
1         loss.backward()
5
1         optimizer.step()
6

```

## 测试

```

1  def greedy_decoder(model, enc_input, start_symbol):
2      """
3          For simplicity, a Greedy Decoder is Beam search when K=1. This is n
4          ecessary for inference as we don't know the
5          target sequence input. Therefore we try to generate the target inpu
6          t word by word, then feed it into the transformer.
7          Starting Reference: http://nlp.seas.harvard.edu/2018/04/03/attention.html#greedy-decoding
8          :param model: Transformer Model
9          :param enc_input: The encoder input
10         :param start_symbol: The start symbol. In this example it is 'S' wh
11         ich corresponds to index 4
12         :return: The target input
13         """
14
15         enc_outputs, enc_self_attns = model.encoder(enc_input)
16
17         dec_input = torch.zeros(1, 0).type_as(enc_input.data)
18
19         terminal = False
20
21         next_symbol = start_symbol
22
23         while not terminal:
24
25             dec_input = torch.cat([dec_input.detach(), torch.tensor([[next_s
26             ymbol]]), dtype=enc_input.dtype).cuda()], -1)
27             dec_outputs, _, _ = model.decoder(dec_input, enc_input, enc_out
28             puts)
29             projected = model.projection(dec_outputs)
30
31             prob = projected.squeeze(0).max(dim=-1, keepdim=False)[1]
32
33             next_word = prob.data[-1]
34
35             next_symbol = next_word
36
37             if next_symbol == tgt_vocab["."]:
38
39                 terminal = True
40

```

```
2         print(next_word)
4
2         return dec_input
5
2
6
2     # Test
7
2     enc_inputs, _, _ = next(iter(loader))
8
2     enc_inputs = enc_inputs.cuda()
9
3     for i in range(len(enc_inputs)):
0
3         greedy_dec_input = greedy_decoder(model, enc_inputs[i].view(1, -1),
1         start_symbol=tgt_vocab["S"])
3         predict, _, _ = model(enc_inputs[i].view(1, -1), greedy_dec_inpu
2         t)
3         predict = predict.data.max(1, keepdim=True)[1]
3
3         print(enc_inputs[i], '->', [idx2word[n.item()] for n in predict.squ
4         eeze()])
```

测试部分代码由网友 qq2014 提供，在此表示感谢

最后给出[完整代码链接](#)（需要科学的力量）

Github 项目地址: [nlp-tutorial](#)

最后编辑于: 2022 年 03 月 20 日

返回文章列表

打赏