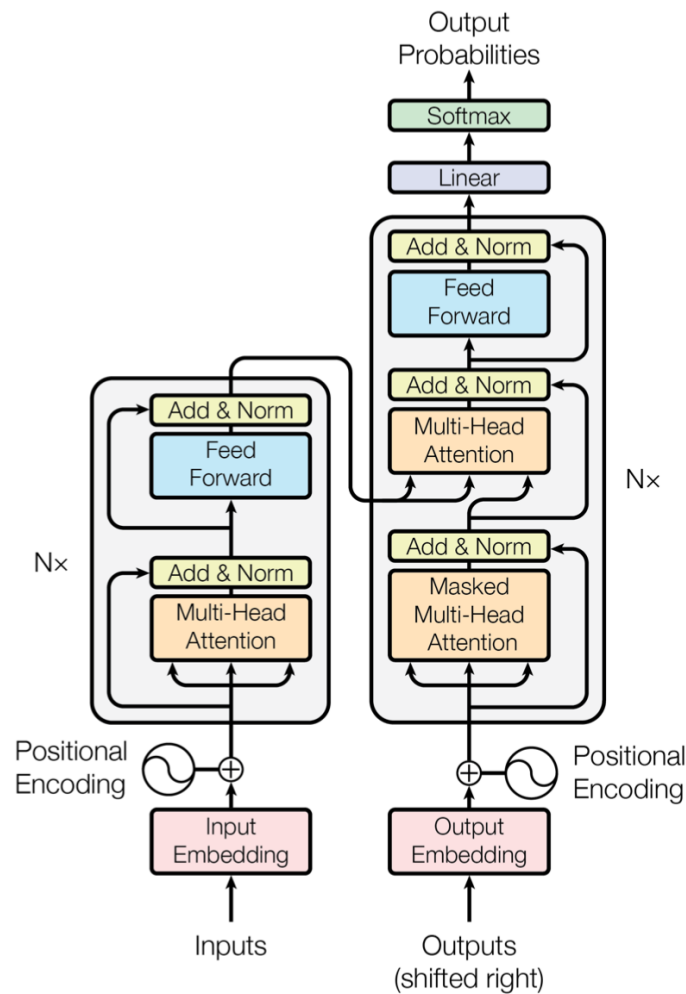


# Table of Contents

- 1 输入序列、目标序列与输出序列
- ▼ 2 词嵌入与位置编码
  - 2.1 输入序列词嵌入与位置编码
  - 2.2 目标序列词嵌入与位置编码
- ▼ 3 编码器Encoder
  - 3.1 编码器结构
  - 3.2 多头注意力机制与缩放点积
  - 3.3 编码器pad掩码
  - 3.4 前馈神经网络
- ▼ 4 解码器Decoder
  - 4.1 解码器结构
  - 4.2 解码器pad掩码、解码器sequence掩码和编码器解码器pad掩码
- ▼ 5 代码实现
  - 5.1 数据处理
  - 5.2 Transformer定义
  - 5.3 编码器定义
  - 5.4 解码器定义
  - 5.5 缩放点积多头注意力机制与前馈神经网络
  - 5.6 掩码与位置编码
  - 5.7 模型训练与验证

## Transformer通用特征提取器



## 1 输入序列、目标序列与输出序列

输入序列  $inputs = (i_1, i_2, \dots, i_p, \dots, i_N)$ ，其中  $i_p \in \mathbb{N}^*$  为输入符号表中的序号。  
 目标序列  $targets = (t_1, t_2, \dots, t_q, \dots, t_M)$ ，其中  $t_q \in \mathbb{N}^*$  为目标符号表中的序号。

$$outputs\_probabilities = Transformer(inputs, targets)$$

其中， $outputs\_probabilities = (o_1, o_2, \dots, o_q, \dots, o_M)$  为预测序列， $o_q \in \mathbb{N}^*$  为目标符号表中的序号。

在自然语言处理任务中，当输入序列与目标序列中的元素较多，通常以句子为单位划分为若干个对应的“输入-目标”子序列进行学习。

## 2 词嵌入与位置编码

### 2.1 输入序列词嵌入与位置编码

输入序列词嵌入  $Embedding(inputs) \in \mathbb{R}^{N \times d_{model}}$ ，其中， $N$  为输入序列长度， $d_{model}$  为词嵌入维度。

输入序列位置编码  $Pos\_Enc(inputs\_position) \in \mathbb{R}^{N \times d_{model}}$ ，

其中， $inputs\_position = (1, 2, \dots, p, \dots, N)$  为输入序列中输入符号对应的位置序号；

$$Pos\_Enc_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$Pos\_Enc_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

其中， $pos \in inputs\_position, i \in (0, 1, \dots, d_{model}/2)$ 。

### 2.2 目标序列词嵌入与位置编码

目标序列词嵌入  $Embedding(targets) \in \mathbb{R}^{M \times d_{model}}$ ，其中  $M$  为目标序列长度， $d_{model}$  为词嵌入维度。

目标序列位置编码  $Pos\_Enc(targets\_position) \in \mathbb{R}^{M \times d_{model}}$ ，

其中， $targets\_position = (1, 2, \dots, q, \dots, M)$  为目标序列的位置序号。

## 3 编码器Encoder

### 3.1 编码器结构

编码器结构：

$$e_0 = Embedding(inputs) + Pos\_Enc(inputs\_position)$$

$$e_l = EncoderLayer(e_{l-1}), l \in [1, n]$$

其中， $e_0 \in \mathbb{R}^{N \times d_{model}}$  为编码器输入， $EncoderLayer(\cdot)$  为编码器层， $n$  为层数， $e_l \in \mathbb{R}^{N \times d_{model}}$  为第  $l$  层编码器层输出。

编码器层  $EncoderLayer$ ：

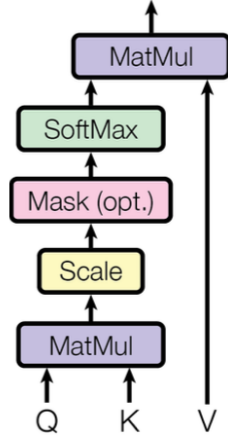
$$e_{mid} = LayerNorm(e_{in} + MultiHeadAttention(e_{in}))$$

$$e_{out} = LayerNorm(e_{mid} + FFN(e_{mid}))$$

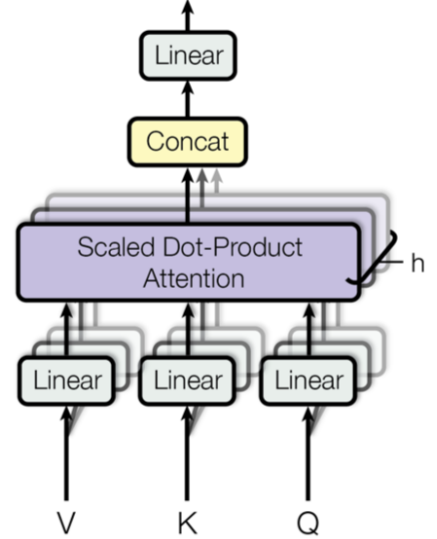
其中， $e_{in} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输入， $e_{out} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输出， $MultiHeadAttention(\cdot)$  为多头注意力机制， $FFN(\cdot)$  为前馈神经网络， $LayerNorm(\cdot)$  为层归一化。

### 3.2 多头注意力机制与缩放点积

## Scaled Dot-Product Attention



## Multi-Head Attention



输入向量序列  $e_{in} = (e_{in1}, e_{in2}, \dots, e_{inN}) \in \mathbb{R}^{N \times d_{model}}$ , 分别得到查询向量序列  $Q = e_{in}$ , 键向量序列  $K = e_{in}$ , 值向量序列  $V = e_{in}$ 。

多头注意力机制

$$MultiHeadAttention(e_{in}) = MultiHead(Q, K, V) = Concat(head_1, \dots, head_h) W^O$$

其中, 多头输出  $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$ , 可学习的参数矩阵

$$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}, W_i^K \in \mathbb{R}^{d_{model} \times d_k}, W_i^V \in \mathbb{R}^{d_{model} \times d_v}, W^O \in \mathbb{R}^{hd_v \times d_{model}}$$

使用缩放点积作为打分函数的自注意力机制

$$Attention(QW_i^Q, KW_i^K, VW_i^V) = softmax\left(\frac{QW_i^Q(KW_i^K)^\top}{\sqrt{d_k}}\right) VW_i^V$$

### 3.3 编码器pad掩码

$$enc\_pad\_mask_j = (e_{j1}, e_{j2}, \dots, e_{jp}, \dots, e_{jN})$$

其中,

$$e_{jp} = \begin{cases} True, & i_p = 0 \\ False, & i_p \neq 0 \end{cases} \quad j = 1, 2, \dots, N$$

$enc\_pad\_mask \in \mathbb{R}^{N \times N}$ ,  $i_p$  为输入序列  $inputs$  对应位置序号。

### 3.4 前馈神经网络

$$FFN(e_{mid}) = ReLU(e_{mid}W_1 + b_1)W_2 + b_2$$

$$= \max(0, e_{mid}W_1 + b_1)W_2 + b_2$$

其中, 参数矩阵  $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}$ ,  $W_2 \in \mathbb{R}^{d_{ff} \times d_{model}}$ , 偏置  $b_1 \in \mathbb{R}^{d_{ff}}$ ,  $b_2 \in \mathbb{R}^{d_{model}}$ 。

## 4 解码器Decoder

### 4.1 解码器结构

$$d_0 = Embedding(targets) + Pos\_Enc(targets\_position)$$

$$d_l = DecoderLayer(d_{l-1}), l \in [1, n]$$

$$outputs\_probabilities = softmax(d_n W)$$

其中,  $d_0 \in \mathbb{R}^{M \times d_{model}}$  为解码器输入,  $DecoderLayer(\cdot)$  为解码器层,  $n$  为层数,  $d_l \in \mathbb{R}^{M \times d_{model}}$  为第  $l$  层解码器层输出,  $W \in \mathbb{R}^{M \times vocab\_size}$  输入输出参数矩阵,  $softmax(\cdot)$  为 softmax 层。

解码器层DecoderLayer:

$$d_{mid1} = LayerNorm(d_{in} + MaskedMultiHeadAttention(d_{in}))$$

$$d_{mid2} = LayerNorm(d_{mid1} + MultiHeadAttention(d_{mid1}, e_{out}))$$

$$d_{out} = LayerNorm(d_{mid2} + FFN(d_{mid2}))$$

其中,  $d_{in} \in \mathbb{R}^{M \times d_{model}}$  为解码器层输入,  $d_{out} \in \mathbb{R}^{M \times d_{model}}$  为解码器层输出,  $MultiHeadAttention(\cdot)$  为多头注意力机制,  $FFN(\cdot)$  为前馈神经网络,  $LayerNorm(\cdot)$  为层归一化。

## 4.2 解码器pad掩码、解码器sequence掩码和编码器解码器pad掩码

解码器pad掩码

$$dec\_pad\_mask_j = (d_{j1}, d_{j2}, \dots, d_{jq}, \dots, d_{jM})$$

其中,

$$d_{jq} = \begin{cases} True, & t_q = 0 \\ False, & t_q \neq 0 \end{cases} \quad j = 1, 2, \dots, M$$

$dec\_pad\_mask \in \mathbb{R}^{M \times M}$ ,  $t_q$  为目标序列  $targets$  对应位置序号。

解码器sequence掩码

$$dec\_sequence\_mask_j = (s_{j1}, s_{j2}, \dots, s_{jl}, \dots, s_{jM})$$

其中,

$$s_{jl} = \begin{cases} 0, & j \geq l \\ 1, & j < l \end{cases} \quad j = 1, 2, \dots, M$$

$dec\_sequence\_mask \in \mathbb{R}^{M \times M}$ , 为非零元素为1的上三角矩阵。

解码器掩码

$$dec\_mask = dec\_pad\_mask + dec\_sequence\_mask$$

编码器解码器pad掩码

$$dec\_enc\_pad\_mask_j = (de_{j1}, de_{j2}, \dots, de_{jp}, \dots, de_{jN})$$

其中,

$$de_{jp} = \begin{cases} True, & i_p = 0 \\ False, & i_p \neq 0 \end{cases} \quad j = 1, 2, \dots, M$$

$dec\_enc\_pad\_mask \in \mathbb{R}^{M \times N}$ ,  $i_p$  为输入序列  $inputs$  对应位置序号。

## 5 代码实现

### 5.1 数据处理

```

In [3]: 1 import numpy as np
        2 import torch
        3 import torch.nn as nn
        4 import torch.optim as optim
        5 from torch.autograd import Variable
        6 import matplotlib.pyplot as plt
        7
        8 dtype = torch.FloatTensor
        9
       10 sentences = ['ich mochte ein bier P', 'S i want a beer', 'i want a beer E']
       11 src_vocab = {'P' : 0, 'ich' : 1, 'mochte' : 2, 'ein' : 3, 'bier' : 4}
       12 src_vocab_size = len(src_vocab)
       13
       14 tgt_vocab = {'P' : 0, 'i' : 1, 'want' : 2, 'a' : 3, 'beer' : 4, 'S' : 5, 'E' : 6}
       15 number_dict = {i: w for i, w in enumerate(tgt_vocab)}
       16 tgt_vocab_size = len(tgt_vocab)
       17
       18 src_len = 5
       19 tgt_len = 5
       20
       21 d_model = 512
       22 d_ff = 2048
       23 d_k = d_v = 64
       24 n_layers = 6
       25 n_heads = 8
       26
       27 def make_batch(sentences):
       28     input_batch = [[src_vocab[w] for w in sentences[0].split()]]
       29     output_batch = [[tgt_vocab[w] for w in sentences[1].split()]]
       30     target_batch = [[tgt_vocab[w] for w in sentences[2].split()]]
       31     return Variable(torch.LongTensor(input_batch)), Variable(torch.LongTensor(output_batch))
       32
executed in 9ms, finished 17:19:54 2020-04-23

```

```

In [4]: 1 make_batch(sentences)
executed in 7ms, finished 17:19:54 2020-04-23

```

```

Out[4]: (tensor([[1, 2, 3, 4, 0]]),
         tensor([[5, 1, 2, 3, 4]]),
         tensor([[1, 2, 3, 4, 6]]))

```

## ▼ 5.2 Transformer定义

```

In [5]: 1 class Transformer(nn.Module):
        2     def __init__(self):
        3         super(Transformer, self).__init__()
        4         self.encoder = Encoder()
        5         self.decoder = Decoder()
        6         self.projection = nn.Linear(d_model, tgt_vocab_size, bias=False)
        7     def forward(self, enc_inputs, dec_inputs):
        8         enc_outputs, enc_self_attns = self.encoder(enc_inputs)
        9         dec_outputs, dec_self_attns, dec_enc_attns = self.decoder(dec_inputs, enc_inputs,
       10         dec_logits = self.projection(dec_outputs) # dec_logits : [batch_size x src_vocab_size]
       11         return dec_logits.view(-1, dec_logits.size(-1)), enc_self_attns, dec_self_attns, dec_enc_attns
       12
executed in 7ms, finished 17:19:56 2020-04-23

```

## ▼ 5.3 编码器定义

```

In [6]: 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 = nn.Embedding.from_pretrained(get_sinusoid_encoding_table(src_len+1,
6             self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])
7
8         def forward(self, enc_inputs): # enc_inputs : [batch_size x source_len]
9             enc_outputs = self.src_emb(enc_inputs) + self.pos_emb(torch.LongTensor([[1,2,3,4,0
10            enc_self_attn_mask = get_attn_pad_mask(enc_inputs, enc_inputs)
11            enc_self_attns = []
12            for layer in self.layers:
13                enc_outputs, enc_self_attn = layer(enc_outputs, enc_self_attn_mask)
14                enc_self_attns.append(enc_self_attn)
15            return enc_outputs, enc_self_attns
16
17 class EncoderLayer(nn.Module):
18     def __init__(self):
19         super(EncoderLayer, self).__init__()
20         self.enc_self_attn = MultiHeadAttention()
21         self.pos_ffn = PoswiseFeedForwardNet()
22
23     def forward(self, enc_inputs, enc_self_attn_mask):
24         enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs, enc_inputs, enc_sel
25         enc_outputs = self.pos_ffn(enc_outputs) # enc_outputs: [batch_size x len_q x d_model]
26         return enc_outputs, attn
27
executed in 9ms, finished 17:19:58 2020-04-23

```

## 5.4 解码器定义

```

In [7]: 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 = nn.Embedding.from_pretrained(get_sinusoid_encoding_table(tgt_len+1,
6             self.layers = nn.ModuleList([DecoderLayer() for _ in range(n_layers)])
7
8         def forward(self, dec_inputs, enc_inputs, enc_outputs): # dec_inputs : [batch_size x t
9             dec_outputs = self.tgt_emb(dec_inputs) + self.pos_emb(torch.LongTensor([[5,1,2,3,4
10            dec_self_attn_pad_mask = get_attn_pad_mask(dec_inputs, dec_inputs)
11            dec_self_attn_subsequent_mask = get_attn_subsequent_mask(dec_inputs)
12            dec_self_attn_mask = torch.gt((dec_self_attn_pad_mask + dec_self_attn_subsequent_m
13
14            dec_enc_attn_mask = get_attn_pad_mask(dec_inputs, enc_inputs)
15
16            dec_self_attns, dec_enc_attns = [], []
17            for layer in self.layers:
18                dec_outputs, dec_self_attn, dec_enc_attn = layer(dec_outputs, enc_outputs, dec
19                dec_self_attns.append(dec_self_attn)
20                dec_enc_attns.append(dec_enc_attn)
21            return dec_outputs, dec_self_attns, dec_enc_attns
22
23 class DecoderLayer(nn.Module):
24     def __init__(self):
25         super(DecoderLayer, self).__init__()
26         self.dec_self_attn = MultiHeadAttention()
27         self.dec_enc_attn = MultiHeadAttention()
28         self.pos_ffn = PoswiseFeedForwardNet()
29
30     def forward(self, dec_inputs, enc_outputs, dec_self_attn_mask, dec_enc_attn_mask):
31         dec_outputs, dec_self_attn = self.dec_self_attn(dec_inputs, dec_inputs, dec_inputs
32         dec_outputs, dec_enc_attn = self.dec_enc_attn(dec_outputs, enc_outputs, enc_output
33         dec_outputs = self.pos_ffn(dec_outputs)
34         return dec_outputs, dec_self_attn, dec_enc_attn
35
executed in 11ms, finished 17:20:00 2020-04-23

```

## 5.5 缩放点积多头注意力机制与前馈神经网络

```

In [8]: 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)
5             self.W_K = nn.Linear(d_model, d_k * n_heads)
6             self.W_V = nn.Linear(d_model, d_v * n_heads)
7
8         def forward(self, Q, K, V, attn_mask):
9             # q: [batch_size x len_q x d_model], k: [batch_size x len_k x d_model], v: [batch_size x len_v x d_model]
10            residual, batch_size = Q, Q.size(0)
11            # (B, S, D) -proj-> (B, S, H, W) -split-> (B, S, H, W) -trans-> (B, H, S, W)
12            q_s = self.W_Q(Q).view(batch_size, -1, n_heads, d_k).transpose(1,2) # q_s: [batch_size x n_heads x len_q x d_k]
13            k_s = self.W_K(K).view(batch_size, -1, n_heads, d_k).transpose(1,2) # k_s: [batch_size x n_heads x len_k x d_k]
14            v_s = self.W_V(V).view(batch_size, -1, n_heads, d_v).transpose(1,2) # v_s: [batch_size x n_heads x len_v x d_v]
15
16            attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1) # attn_mask : [batch_size x n_heads x len_q x len_k]
17
18            # context: [batch_size x n_heads x len_q x d_v], attn: [batch_size x n_heads x len_q x len_k]
19            context, attn = ScaledDotProductAttention()(q_s, k_s, v_s, attn_mask)
20            context = context.transpose(1, 2).contiguous().view(batch_size, -1, n_heads * d_v)
21            output = nn.Linear(n_heads * d_v, d_model)(context)
22            return nn.LayerNorm(d_model)(output + residual), attn # output: [batch_size x len_q x d_model], attn: [batch_size x len_q x len_k]
23
24 class ScaledDotProductAttention(nn.Module):
25     def __init__(self):
26         super(ScaledDotProductAttention, self).__init__()
27
28     def forward(self, Q, K, V, attn_mask):
29         scores = torch.matmul(Q, K.transpose(-1, -2)) / np.sqrt(d_k) # scores : [batch_size x n_heads x len_q x len_k]
30         scores.masked_fill_(attn_mask, -1e9) # Fills elements of self tensor with value where mask == 1
31         attn = nn.Softmax(dim=-1)(scores)
32         context = torch.matmul(attn, V)
33         return context, attn
34
35 class PoswiseFeedForwardNet(nn.Module):
36     def __init__(self):
37         super(PoswiseFeedForwardNet, self).__init__()
38         self.conv1 = nn.Conv1d(in_channels=d_model, out_channels=d_ff, kernel_size=1)
39         self.conv2 = nn.Conv1d(in_channels=d_ff, out_channels=d_model, kernel_size=1)
40
41     def forward(self, inputs):
42         residual = inputs # inputs : [batch_size, len_q, d_model]
43         output = nn.ReLU()(self.conv1(inputs.transpose(1, 2)))
44         output = self.conv2(output).transpose(1, 2)
45         return nn.LayerNorm(d_model)(output + residual)
46

```

executed in 15ms, finished 17:20:02 2020-04-23

## 5.6 掩码与位置编码

```
In [9]: 1 def get_sinusoid_encoding_table(n_position, d_model):
2         def cal_angle(position, hid_idx):
3             return position / np.power(10000, 2 * (hid_idx // 2) / d_model)
4         def get_posi_angle_vec(position):
5             return [cal_angle(position, hid_j) for hid_j in range(d_model)]
6
7         sinusoid_table = np.array([get_posi_angle_vec(pos_i) for pos_i in range(n_position)])
8         sinusoid_table[:, 0::2] = np.sin(sinusoid_table[:, 0::2]) # dim 2i
9         sinusoid_table[:, 1::2] = np.cos(sinusoid_table[:, 1::2]) # dim 2i+1
10        return torch.FloatTensor(sinusoid_table)
11
12    def get_attn_pad_mask(seq_q, seq_k):
13        batch_size, len_q = seq_q.size()
14        batch_size, len_k = seq_k.size()
15        # eq(zero) is PAD token
16        pad_attn_mask = seq_k.data.eq(0).unsqueeze(1) # batch_size x 1 x len_k(=len_q), one is 1
17        return pad_attn_mask.expand(batch_size, len_q, len_k) # batch_size x len_q x len_k
18
19    def get_attn_subsequent_mask(seq):
20        attn_shape = [seq.size(0), seq.size(1), seq.size(1)]
21        subsequent_mask = np.triu(np.ones(attn_shape), k=1)
22        subsequent_mask = torch.from_numpy(subsequent_mask).byte()
23        return subsequent_mask
24
```

executed in 11ms, finished 17:20:03 2020-04-23

## 5.7 模型训练与验证

```
In [10]: 1 model = Transformer()
2
3 criterion = nn.CrossEntropyLoss()
4 optimizer = optim.Adam(model.parameters(), lr=0.001)
5
6 for epoch in range(20):
7     optimizer.zero_grad()
8     enc_inputs, dec_inputs, target_batch = make_batch(sentences)
9     outputs, enc_self_attns, dec_self_attns, dec_enc_attns = model(enc_inputs, dec_inputs)
10    loss = criterion(outputs, target_batch.contiguous().view(-1))
11    print('Epoch:', '%04d' % (epoch + 1), 'cost =', '{:.6f}'.format(loss))
12    loss.backward()
13    optimizer.step()
14
```

executed in 7.39s, finished 17:20:12 2020-04-23

```
Epoch: 0001 cost = 1.946148
Epoch: 0002 cost = 0.061065
Epoch: 0003 cost = 0.061755
Epoch: 0004 cost = 0.040034
Epoch: 0005 cost = 0.030705
Epoch: 0006 cost = 0.001976
Epoch: 0007 cost = 0.001366
Epoch: 0008 cost = 0.001905
Epoch: 0009 cost = 0.003293
Epoch: 0010 cost = 0.004082
Epoch: 0011 cost = 0.000338
Epoch: 0012 cost = 0.002591
Epoch: 0013 cost = 0.005281
Epoch: 0014 cost = 0.000948
Epoch: 0015 cost = 0.000256
Epoch: 0016 cost = 0.000409
Epoch: 0017 cost = 0.000859
Epoch: 0018 cost = 0.024325
Epoch: 0019 cost = 0.007099
Epoch: 0020 cost = 0.020045
```



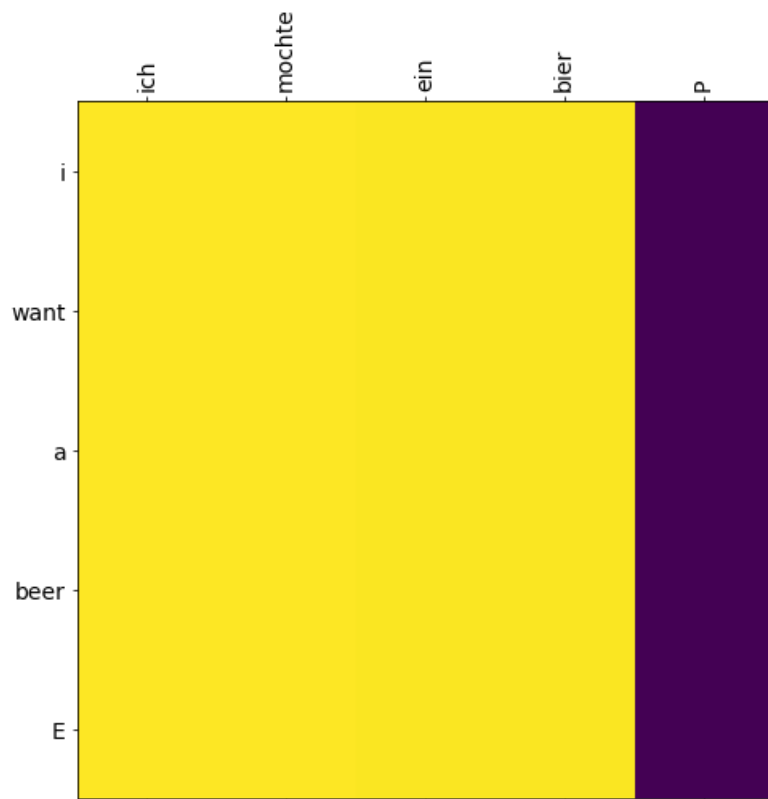
```

In [11]: 1 def showgraph(attn):
2         attn = attn[-1].squeeze(0)[0]
3         attn = attn.squeeze(0).data.numpy()
4         fig = plt.figure(figsize=(n_heads, n_heads)) # [n_heads, n_heads]
5         ax = fig.add_subplot(1, 1, 1)
6         ax.matshow(attn, cmap='viridis')
7         ax.set_xticklabels(['']+sentences[0].split(), fontdict={'fontsize': 14}, rotation=90)
8         ax.set_yticklabels(['']+sentences[2].split(), fontdict={'fontsize': 14})
9         plt.show()
10
11 predict, _, _, _ = model(enc_inputs, dec_inputs)
12
13 predict = predict.data.max(1, keepdim=True)[1]
14 print(sentences[0], '->', [number_dict[n.item()] for n in predict.squeeze()])
15
16 print('first head of last state enc_self_attns')
17 showgraph(enc_self_attns)
18
19 print('first head of last state dec_self_attns')
20 showgraph(dec_self_attns)
21
22 print('first head of last state dec_enc_attns')
23 showgraph(dec_enc_attns)
24

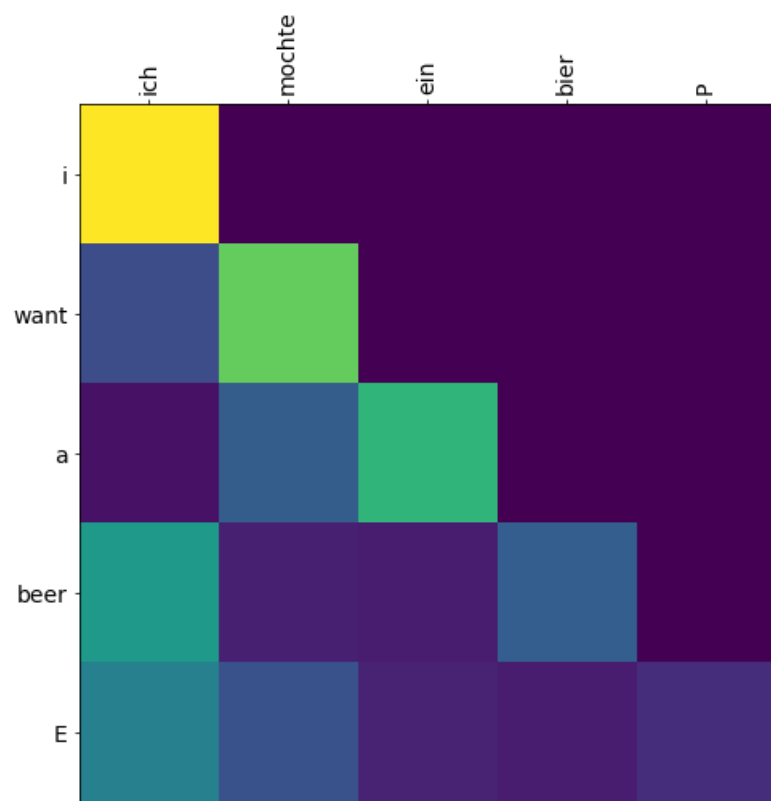
```

executed in 635ms, finished 17:20:13 2020-04-23

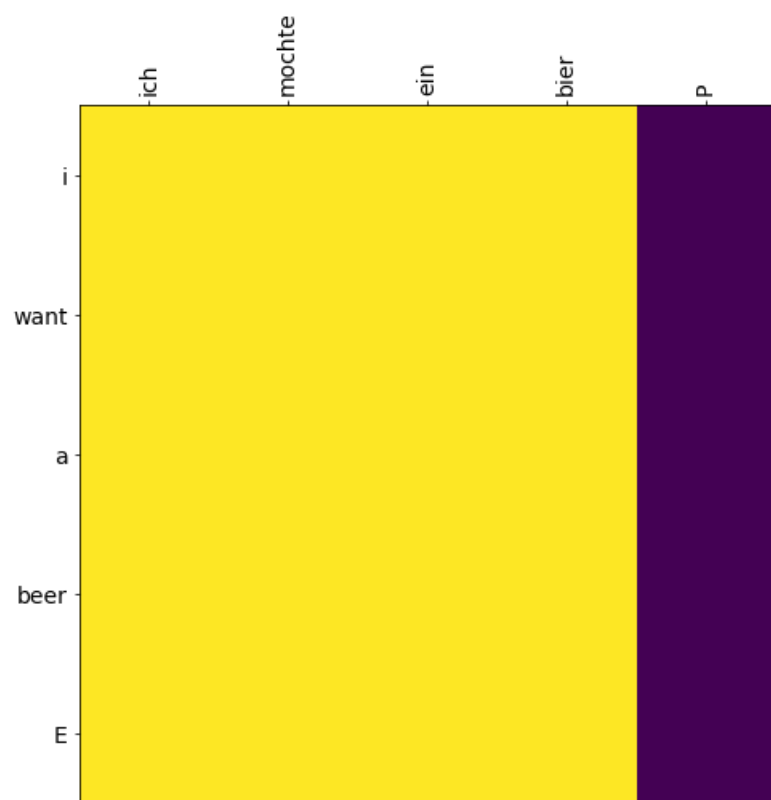
ich mochte ein bier P -> ['i', 'want', 'a', 'beer', 'E']  
first head of last state enc\_self\_attns



first head of last state dec\_self\_attns



first head of last state dec\_enc\_attns



In [ ]:

1

