

如何构建一个基本的Transformer模型 (英中翻译)

发表评论 / ChatGPT, GPT, OpenAI

一个简化版的Transformer模型训练程序示例，使用PyTorch框架。这个示例展示了如何构建一个基本的Transformer模型，用于一个简单的序列到序列的任务（例如，机器翻译或文本生成）。注意，这个例子是为了演示目的而简化的，实际应用中可能需要更复杂的数据处理、模型架构调整和训练策略。

0. 准备环境

确保安装了PyTorch。可以通过`pip install torch`安装。

```
1 <code>pip install torch</code>
```

1.

1. 准备训练样本

我们将定义一小批简单的英文句子及其对应的中文翻译。这里使用的是极其简化的数据集，仅用于演示目的。

```
1 # 英文到中文的简单句子对
2 english_sentences = &#91;
3     "Hello, how are you?", # 你好，你怎么样？
4     "I am learning translation.", # 我在学习翻译。
5     "This is a pen.", # 这是一支笔。
6     "What is your name?", # 你叫什么名字？
7     "I love programming." # 我爱编程。
8 ]
9
10 chinese_sentences = &#91;
11     "你好，你怎么样？",
12     "我在学习翻译。",
13     "这是一支笔。",
14     "你叫什么名字？",
15     "我爱编程。"
16 ]
17
18 # 假设我们已经有了英文和中文的词汇表和编码函数（在实际应用中，需要根据实际数据构建）
19 # 为了简化，我们这里不实现这一部分，而是直接使用英文和中文句子的索引表示
20
21
```

2. 训练代码

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import math
5 import jieba
6 import pickle
7
8 class TransformerModel(nn.Module):
9     def __init__(self, ntoken, ninp, nhead, nhid, nlayers, dropout=0.5):
```

```

10     super(TransformerModel, self).__init__()
11     self.model_type = 'Transformer'
12     self.pos_encoder = PositionalEncoding(ninp, dropout)
13     encoder_layers = nn.TransformerEncoderLayer(d_model=ninp, nhead=nhead, dim_feedforward=dim_feedforward)
14     self.transformer_encoder = nn.TransformerEncoder(encoder_layers, num_layers=nlayers)
15     self.encoder = nn.Embedding(ntoken, ninp)
16     self.ninp = ninp
17     self.decoder = nn.Linear(ninp, ntoken)
18     self.init_weights()
19
20     def init_weights(self):
21         initrange = 0.1
22         self.encoder.weight.data.uniform_(-initrange, initrange)
23         self.decoder.bias.data.zero_()
24         self.decoder.weight.data.uniform_(-initrange, initrange)
25
26     def forward(self, src):
27         src = self.encoder(src) * math.sqrt(self.ninp)
28         src = self.pos_encoder(src)
29         output = self.transformer_encoder(src)
30         output = self.decoder(output)
31         return output
32
33 class PositionalEncoding(nn.Module):
34     def __init__(self, d_model, dropout=0.1, max_len=5000):
35         super(PositionalEncoding, self).__init__()
36         self.dropout = nn.Dropout(p=dropout)
37
38         position = torch.arange(max_len).unsqueeze(1)
39         div_term = torch.exp(torch.arange(0, d_model, 2) * -(math.log(10000.0) / d_model))
40         pe = torch.zeros(max_len, d_model)
41         pe[:, 0::2] = torch.sin(position * div_term)
42         pe[:, 1::2] = torch.cos(position * div_term)
43         self.register_buffer('pe', pe.unsqueeze(0))
44
45     def forward(self, x):
46         x = x + self.pe[:, :x.size(1)]
47         return self.dropout(x)
48
49 # 英文到中文的简单句子对
50 english_sentences = &#91;
51     "Hello, how are you?", # 你好,你怎么样?
52     "I am learning translation.", # 我在学习翻译。
53     "This is a pen.", # 这是一支笔。
54     "What is your name?", # 你叫什么名字?
55     "I love programming." # 我爱编程。
56 ]
57
58 chinese_sentences = &#91;
59     "你好,你怎么样?",
60     "我在学习翻译。",
61     "这是一支笔。",
62     "你叫什么名字?",
63     "我爱编程。"
64 ]
65
66 # 定义结束标记
67 end_token = '<eos>'
68
69 # 构建词汇表的函数
70 def build_english_vocab(sentences):
71     vocab = set(word for sentence in sentences for word in sentence.split())
72     # 添加结束标记到词汇表中
73     vocab.add(end_token)
74     return {word: i for i, word in enumerate(vocab)}

```

```

75
76 # 使用jieba进行分词
77 def build_chinese_vocab(sentences):
78     vocab = set(word for sentence in sentences for word in jieba.cut(sentence))
79     # 添加结束标记到词汇表中
80     vocab.add(end_token)
81     return {word: i for i, word in enumerate(vocab)}
82
83 def encode_english(sentence, vocab, max_len):
84     words = sentence.split()[:max_len - 1] # 保留一个位置给结束标记
85     words.append(end_token) # 添加结束标记
86     return ['<vocab.get(word, vocab['<unk>']) for word in words]
87
88 def encode_chinese(sentence, vocab, max_len):
89     words = list(jieba.cut(sentence))[:max_len - 1] # 保留一个位置给结束标记
90     words.append(end_token) # 添加结束标记
91     return ['<vocab.get(word, vocab['<unk>']) for word in words]
92
93 # 参数设置
94 ntokens = 1000 # 词汇表大小
95 emsize = 200 # 嵌入维度
96 nhid = 200 # 前馈网络的维度
97 nlayers = 2 # Transformer层的数量
98 nhead = 2 # 多头注意力的头数
99 dropout = 0.2 # dropout的比例
100
101 english_vocab = build_english_vocab(english_sentences)
102 chinese_vocab = build_chinese_vocab(chinese_sentences)
103
104 # 编码函数
105 # 假设词汇表中包含<unk>
106 english_vocab['<unk>'] = len(english_vocab)
107 chinese_vocab['<unk>'] = len(chinese_vocab)
108
109 # 定义最大序列长度
110 max_seq_length = max(len(sentence.split()) for sentence in english_sentences)
111
112 # 编码英文和法文句子
113 encoded_english_sentences = ['<encode_english(sentence, english_vocab, max_seq_length) for
114 encoded_chinese_sentences = ['<encode_chinese(sentence, chinese_vocab, max_seq_length) for
115
116 # 词汇表大小
117 english_vocab_size = len(english_vocab)
118 chinese_vocab_size = len(chinese_vocab)
119
120 print(f"english_vocab_size {english_vocab_size}")
121 print(f"chinese_vocab_size {chinese_vocab_size}")
122
123 ntokens = chinese_vocab_size # 将中文词汇表的大小用作 ntokens
124
125 model = TransformerModel(ntokens, emsize, nhead, nhid, nlayers, dropout)
126
127 # 训练代码 (伪代码)
128 criterion = nn.CrossEntropyLoss()
129 optimizer = optim.SGD(model.parameters(), lr=0.01)
130
131 # 假设已经定义了模型、损失函数和优化器
132 # 注意：在实际应用中，你需要根据任务调整模型的输入输出维度以及其他参数
133
134 epochs = 2000
135
136 for epoch in range(epochs):
137     total_loss = 0
138     for eng, chi in zip(encoded_english_sentences, encoded_chinese_sentences):

```

```

140     model.train()
141     optimizer.zero_grad()
142
143     # 调整输入序列的形状以匹配 batch_first=True
144     # 批次大小, 序列长度]
145     src = torch.tensor(eng, dtype=torch.long) # 添加额外的维度来表示批次大小
146     tgt = torch.tensor(chi, dtype=torch.long)
147
148     output = model(src)
149
150     # 重塑输出和目标以适应交叉熵损失
151     output_resaped = output.view(-1, ntokens)
152     tgt_resaped = tgt.view(-1)
153
154     loss = criterion(output_resaped, tgt_resaped)
155     loss.backward()
156     optimizer.step()
157
158     total_loss += loss.item()
159
160     if epoch % 100 == 0:
161         print(f"Epoch {epoch+1}, Loss: {total_loss / len(encoded_english_sentences)}")
162
163     # 保存模型的状态字典
164     model_path = f"models/gpt_model_1_{epochs}.pth"
165     torch.save(model.state_dict(), model_path)
166     print(f"Model saved to {model_path}")
167
168     # 保存 ntokens
169     print(f"ntokens {ntokens}")
170     with open(f"models/ntokens_1_{epochs}.pkl", "wb") as f:
171         pickle.dump(ntokens, f)
172
173     # 保存 max_seq_length
174     print(f"max_seq_length {max_seq_length}")
175     with open(f"models/max_seq_length_1_{epochs}.pkl", "wb") as f:
176         pickle.dump(max_seq_length, f)
177
178     # 保存英文词汇表
179     print(f"english_vocab {english_vocab}")
180     with open(f"models/english_vocab_1_{epochs}.pkl", "wb") as f:
181         pickle.dump(english_vocab, f)
182
183     # 保存中文词汇表
184     print(f"chinese_vocab {chinese_vocab}")
185     with open(f"models/chinese_vocab_1_{epochs}.pkl", "wb") as f:
186         pickle.dump(chinese_vocab, f)
187
188
189     # 输出模型参数数量
190     print(sum(p.numel() for p in model.parameters())/1e6, 'M parameters')
191

```

3. 加载模型产生文本

```

1  import torch
2  import torch.nn.functional as F
3  import torch.nn as nn
4  import math
5  import pickle
6
7  class TransformerModel(nn.Module):
8      def __init__(self, ntoken, ninp, nhead, nhid, nlayers, dropout=0.5):

```

```

9         super(TransformerModel, self).__init__()
10        self.model_type = 'Transformer'
11        self.pos_encoder = PositionalEncoding(ninp, dropout)
12        encoder_layers = nn.TransformerEncoderLayer(d_model=ninp, nhead=nhead, dim_feedforward=dim_feedforward)
13        self.transformer_encoder = nn.TransformerEncoder(encoder_layers, num_layers=nlayers)
14        self.encoder = nn.Embedding(ntoken, ninp)
15        self.ninp = ninp
16        self.decoder = nn.Linear(ninp, ntoken)
17        self.init_weights()
18
19    def init_weights(self):
20        initrange = 0.1
21        self.encoder.weight.data.uniform_(-initrange, initrange)
22        self.decoder.bias.data.zero_()
23        self.decoder.weight.data.uniform_(-initrange, initrange)
24
25    def forward(self, src):
26        src = self.encoder(src) * math.sqrt(self.ninp)
27        src = self.pos_encoder(src)
28        output = self.transformer_encoder(src)
29        output = self.decoder(output)
30        return output
31
32    class PositionalEncoding(nn.Module):
33        def __init__(self, d_model, dropout=0.1, max_len=5000):
34            super(PositionalEncoding, self).__init__()
35            self.dropout = nn.Dropout(p=dropout)
36
37            position = torch.arange(max_len).unsqueeze(1)
38            div_term = torch.exp(torch.arange(0, d_model, 2) * -(math.log(10000.0) / d_model))
39            pe = torch.zeros(max_len, d_model)
40            pe[:, 0::2] = torch.sin(position * div_term)
41            pe[:, 1::2] = torch.cos(position * div_term)
42            self.register_buffer('pe', pe.unsqueeze(0))
43
44        def forward(self, x):
45            x = x + self.pe[:, :x.size(1)]
46            return self.dropout(x)
47
48    end_token = '<eos>'
49
50    def encode_english(sentence, vocab, max_len):
51        words = sentence.split()[:max_len - 1] # 保留一个位置给结束标记
52        words.append(end_token) # 添加结束标记
53        return [vocab.get(word, vocab['<unk>']) for word in words]
54
55    def decode_chinese(indices, vocab):
56        words = [vocab.keys()[vocab.values().index(idx)] for idx in indices]
57        # 去除结束标记之后的部分
58        if end_token in words:
59            words = words[:words.index(end_token)]
60        return ''.join(words)
61
62    # 加载模型
63    ntokens = 1000 # 假设词汇表大小为1000
64    emsize = 200 # 嵌入维度
65    nhid = 200 # 前馈网络的维度
66    nlayers = 2 # Transformer层的数量
67    nhead = 2 # 多头注意力的头数
68    dropout = 0.2 # dropout的比例
69
70    epochs = 2000
71
72    # 加载 ntokens
73    with open(f"models/ntokens_1_{epochs}.pkl", "rb") as f:

```

```

74     ntokens = pickle.load(f)
75     print(f"ntokens {ntokens}")
76
77     # 加载 max_seq_length
78     with open(f"models/max_seq_length_1_{epochs}.pkl", "rb") as f:
79         max_seq_length = pickle.load(f)
80     print(f"max_seq_length {max_seq_length}")
81
82     # 加载英文词汇表
83     with open(f"models/english_vocab_1_{epochs}.pkl", "rb") as f:
84         english_vocab = pickle.load(f)
85     print(f"english_vocab {english_vocab}")
86
87     # 加载中文词汇表
88     with open(f"models/chinese_vocab_1_{epochs}.pkl", "rb") as f:
89         chinese_vocab = pickle.load(f)
90     print(f"chinese_vocab {chinese_vocab}")
91
92     model = TransformerModel(ntokens, emsize, nhead, nhid, nlayers, dropout)
93
94     model_path = f"models/gpt_model_1_{epochs}.pth"
95     model.load_state_dict(torch.load(model_path))
96     model.eval()
97
98     # 定义生成文本的函数
99     def generate_text(model, start_text, max_len=50):
100         with torch.no_grad():
101             tokenized_start_text = encode_english(start_text, english_vocab, max_seq_length)
102             input_seq = torch.tensor(&#91;tokenized_start_text], dtype=torch.long)
103             for _ in range(max_len):
104                 output = model(input_seq)
105                 predicted_token = torch.argmax(output&#91;::, -1, :], dim=-1) # 获取最后一个时间步
106                 input_seq = torch.cat(&#91;input_seq, predicted_token.unsqueeze(0)], dim=1) # 拼接
107                 if predicted_token.item() == end_token: # 如果生成了终止符号, 停止生成
108                     break
109             generated_text = decode_chinese(input_seq.squeeze().tolist(), chinese_vocab)
110             return generated_text
111
112     # 使用生成文本的函数
113     start_text = "Hello, how are you?" # 初始文本
114     generated_text = generate_text(model, start_text)
115     print("Generated text:", generated_text)
116
117

```