

从头开始实现一个简化的版本的GPT模型

发表评论 / ChatGPT, GPT, OpenAI

如果你想从头开始实现一个简化的版本的GPT模型，而不依赖于现成的GPT-2模型库，你可以采用PyTorch这样的深度学习框架。下面是一个非常基础的例子，展示了如何实现一个简化的Transformer模型架构，这是构建GPT模型的基础。

这个例子将不会覆盖GPT-2的所有复杂性和特性，但可以提供一个起点，帮助你理解如何从头开始构建类似GPT的模型。

1. 基础Transformer块

首先，我们定义一个基础的Transformer块，它是构成GPT模型的基本单元。这个块将包括自注意力机制和前馈神经网络。

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torch.optim as optim
5 import re
6 import jieba
7
8 from torch.utils.data import Dataset, DataLoader
9 from torch.nn.utils.rnn import pad_sequence
10
11
12 class SelfAttention(nn.Module):
13     def __init__(self, embed_size, heads):
14         super(SelfAttention, self).__init__()
15         self.embed_size = embed_size
16         self.heads = heads
17         self.head_dim = embed_size // heads
18
19         assert self.head_dim * heads == embed_size, "Embedding size needs to be divisible by"
20
21         self.values = nn.Linear(embed_size, embed_size, bias=False)
22         self.keys = nn.Linear(embed_size, embed_size, bias=False)
23         self.queries = nn.Linear(embed_size, embed_size, bias=False)
24         self.fc_out = nn.Linear(heads * self.head_dim, embed_size)
25
26     def forward(self, value, key, query, mask):
27         N = query.shape[0]
28         value_len, key_len, query_len = value.shape[1], key.shape[1], query.shape[1]
29
30         # Split the embedding into `self.heads` pieces
31         values = self.values(value).view(N, value_len, self.heads, self.head_dim)
32         keys = self.keys(key).view(N, key_len, self.heads, self.head_dim)
33         queries = self.queries(query).view(N, query_len, self.heads, self.head_dim)
34
35         # Transpose for attention dot product: from [N, value_len, self.heads, self.head_dim]
36         # to [N, self.heads, value_len, self.head_dim] to match the shape for `torch.einsum`
37         values = values.transpose(1, 2)
38         keys = keys.transpose(1, 2)
39         queries = queries.transpose(1, 2)
40
41         # Attention mechanism
42         energy = torch.einsum("nqhd,nkhd->nhqk", [queries, keys])
43
44         if mask is not None:
```

```

45         energy = energy.masked_fill(mask == 0, float("-1e20"))
46
47         attention = torch.softmax(energy / (self.embed_size ** (1 / 2)), dim=3)
48
49         # 重塑前进行张量乘法, 然后重塑回 &#91;batch_size, seq_len, heads * head_dim]
50         out = torch.einsum("nhql,nlhd->nqhd", &#91;attention, values]).reshape(
51             N, query_len, self.heads * self.head_dim
52         )
53
54         out = self.fc_out(out)
55         return out
56
57     class TransformerBlock(nn.Module):
58         def __init__(self, embed_size, heads, dropout, forward_expansion):
59             super(TransformerBlock, self).__init__()
60             self.attention = SelfAttention(embed_size, heads)
61             self.norm1 = nn.LayerNorm(embed_size)
62             self.norm2 = nn.LayerNorm(embed_size)
63
64             self.feed_forward = nn.Sequential(
65                 nn.Linear(embed_size, forward_expansion * embed_size),
66                 nn.ReLU(),
67                 nn.Linear(forward_expansion * embed_size, embed_size),
68             )
69
70             self.dropout = nn.Dropout(dropout)
71
72         def forward(self, value, key, query, mask):
73             attention = self.attention(value, key, query, mask)
74
75             # Add skip connection, followed by layer normalization
76             x = self.norm1(attention + query)
77             forward = self.feed_forward(x)
78             out = self.norm2(forward + x) # Add skip connection, followed by layer normalization
79             return out
80

```

2. 简化版的GPT模型

接下来, 我们定义一个简化版的GPT模型, 它利用上面定义的Transformer块。

```

1  class GPT(nn.Module):
2      def __init__(self, embed_size, num_layers, heads, forward_expansion, dropout, vocab_size,
3          super(GPT, self).__init__()
4          self.embed_size = embed_size
5          self.transformer_blocks = nn.ModuleList(
6              &#91;
7                  TransformerBlock(
8                      embed_size,
9                      heads,
10                     dropout=dropout,
11                     forward_expansion=forward_expansion,
12                 )
13                 for _ in range(num_layers)
14             ]
15         )
16
17         self.word_embedding = nn.Embedding(vocab_size, embed_size)
18         self.position_embedding = nn.Embedding(max_length, embed_size)
19
20     def forward(self, x, mask):
21         N, seq_length = x.shape
22         print(f"Input shape: {x.shape}") # 打印输入形状

```

```

23
24     positions = torch.arange(0, seq_length).expand(N, seq_length).to(x.device)
25     out = self.word_embedding(x) + self.position_embedding(positions)
26     print(f"After embedding and position shape: {out.shape}") # 打印嵌入和位置编码后的形状
27
28     for layer in self.transformer_blocks:
29         out = layer(out, out, out, mask)
30         print(f"After transformer block shape: {out.shape}") # 打印经过每个Transformer块后
31
32     return out
33

```

继续前面的简化版GPT模型实现，下面提供一个基本的训练框架。这个例子将展示如何准备数据、定义损失函数、选择优化器，并执行训练循环。请注意，这是一个高度简化的例子，旨在演示基本概念。

3. 准备数据

假设你已经有了一个文本数据集，并且你已经进行了预处理（例如，分词和转换为词汇索引）。为了简单起见，这里不展示数据预处理的代码。我们将直接从创建数据加载器开始。

```

1  def clean_text_mixed_with_symbols(text):
2      # 保留中文、英文字符、数字和常见的标点符号
3      # 注意：根据需要，你可以在这里添加或删除特定的符号
4      text = re.sub(r'&#91;\u4e00-\u9fffA-Za-z0-9,。! ?、; : " ' ( ) 《 》 [ ] _...]+', ' ', text)
5      return text.strip()
6
7  def preprocess_text_mixed_with_symbols(text):
8      text = clean_text_mixed_with_symbols(text)
9      tokens = &#91;]
10     for token in jieba.cut(text, cut_all=False):
11         token = token.strip()
12         if token:
13             tokens.append(token)
14     return tokens
15
16  def load_and_preprocess_data(file_paths):
17      # 这里简化处理，具体实现依据你的需求定
18      texts = &#91;]
19      for file_path in file_paths:
20          with open(file_path, 'r', encoding='utf-8') as file:
21              text = file.read()
22              # 添加文本清洗和预处理逻辑
23              processed_text = preprocess_text_mixed_with_symbols(text)
24              texts.append(processed_text)
25      return texts
26
27
28  class TextDataset(Dataset):
29      def __init__(self, indexed_texts, vocab_size):
30          self.texts = &#91;torch.tensor(text, dtype=torch.long) for text in indexed_texts] #
31          self.vocab_size = vocab_size
32
33      def __len__(self):
34          return len(self.texts)
35
36      def __getitem__(self, idx):
37          return self.texts[&#91;idx]
38
39      def collate_fn(self, batch):
40          input_ids = &#91;item[&#91;-1] for item in batch]
41          target_ids = &#91;item[&#91;1:] for item in batch]
42          input_ids_padded = pad_sequence(input_ids, batch_first=True, padding_value=0)

```

```

43     target_ids_padded = pad_sequence(target_ids, batch_first=True, padding_value=0)
44     return input_ids_padded, target_ids_padded
45
46
47 def build_vocab(texts):
48     vocab = set(token for text in texts for token in text)
49     vocab_to_index = {word: i for i, word in enumerate(vocab, start=1)} # 从1开始编号
50     return vocab_to_index
51
52
53 def index_text(text, vocab_to_index):
54     return [vocab_to_index[token] for token in text if token in vocab_to_index]
55

```

4. 定义模型、损失函数和优化器

```

1 # 实例化模型
2 model = GPT(
3     embed_size=embed_size,
4     num_layers=num_layers,
5     heads=heads,
6     forward_expansion=forward_expansion,
7     dropout=dropout,
8     vocab_size=vocab_size,
9     max_length=max_length
10 )
11
12 loss_fn = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.0001)
14

```

5. 训练循环

最后，我们执行训练循环，每个批次处理数据，计算损失，并更新模型的权重。

```

1 def train(model, dataloader, loss_fn, optimizer, device, epochs):
2     model.train()
3     model.to(device)
4
5     for epoch in range(epochs):
6         for batch_idx, (input_ids, target_ids) in enumerate(dataloader):
7             input_ids = input_ids.to(device)
8             target_ids = target_ids.to(device)
9
10            # 前向传播
11            predictions = model(input_ids, mask=None) # 这里简化处理，没有使用mask
12            predictions = predictions.view(-1, predictions.size(-1))
13            target_ids = target_ids.view(-1)
14
15            # 计算损失
16            loss = loss_fn(predictions, target_ids)
17
18            # 反向传播和优化
19            optimizer.zero_grad()
20            loss.backward()
21            optimizer.step()
22
23            if batch_idx % 100 == 0:
24                print(f"Epoch {epoch} Batch {batch_idx} Loss {loss.item()}")
25

```

这段代码展示了如何设置和执行模型的训练过程。请注意，这只是一个起点，真实世界的应用可能需要更复杂的数据处理、模型调参、正则化策略、以及训练过程监控。此外，为了处理大规模数据集和模型，可能还需要考虑分布式训练和模型并行化。

6. 完整的训练代码

代码如下：

```

1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F
4  import torch.optim as optim
5  import re
6  import jieba
7
8  from torch.utils.data import Dataset, DataLoader
9  from torch.nn.utils.rnn import pad_sequence
10
11
12  class SelfAttention(nn.Module):
13      def __init__(self, embed_size, heads):
14          super(SelfAttention, self).__init__()
15          self.embed_size = embed_size
16          self.heads = heads
17          self.head_dim = embed_size // heads
18
19          assert self.head_dim * heads == embed_size, "Embedding size needs to be divisible by"
20
21          self.values = nn.Linear(embed_size, embed_size, bias=False)
22          self.keys = nn.Linear(embed_size, embed_size, bias=False)
23          self.queries = nn.Linear(embed_size, embed_size, bias=False)
24          self.fc_out = nn.Linear(heads * self.head_dim, embed_size)
25
26      def forward(self, value, key, query, mask):
27          N = query.shape[0]
28          value_len, key_len, query_len = value.shape[1], key.shape[1], query.shape[1]
29
30          # Split the embedding into `self.heads` pieces
31          values = self.values(value).view(N, value_len, self.heads, self.head_dim)
32          keys = self.keys(key).view(N, key_len, self.heads, self.head_dim)
33          queries = self.queries(query).view(N, query_len, self.heads, self.head_dim)
34
35          # Transpose for attention dot product: from [N, value_len, self.heads, self.head_dim]
36          # to [N, self.heads, value_len, self.head_dim] to match the shape for `torch.einsum`
37          values = values.transpose(1, 2)
38          keys = keys.transpose(1, 2)
39          queries = queries.transpose(1, 2)
40
41          # Attention mechanism
42          energy = torch.einsum("nqhd,nkhd->nhqk", [queries, keys])
43
44          if mask is not None:
45              energy = energy.masked_fill(mask == 0, float("-1e20"))
46
47          attention = torch.softmax(energy / (self.embed_size ** (1 / 2)), dim=3)
48
49          # 重塑前进行张量乘法，然后重塑回 [batch_size, seq_len, heads * head_dim]
50          out = torch.einsum("nhql,nlhd->nqhd", [attention, values]).reshape(
51              N, query_len, self.heads * self.head_dim
52          )
53
54          out = self.fc_out(out)

```

```

55     return out
56
57 class TransformerBlock(nn.Module):
58     def __init__(self, embed_size, heads, dropout, forward_expansion):
59         super(TransformerBlock, self).__init__()
60         self.attention = SelfAttention(embed_size, heads)
61         self.norm1 = nn.LayerNorm(embed_size)
62         self.norm2 = nn.LayerNorm(embed_size)
63
64         self.feed_forward = nn.Sequential(
65             nn.Linear(embed_size, forward_expansion * embed_size),
66             nn.ReLU(),
67             nn.Linear(forward_expansion * embed_size, embed_size),
68         )
69
70         self.dropout = nn.Dropout(dropout)
71
72     def forward(self, value, key, query, mask):
73         attention = self.attention(value, key, query, mask)
74
75         # Add skip connection, followed by layer normalization
76         x = self.norm1(attention + query)
77         forward = self.feed_forward(x)
78         out = self.norm2(forward + x) # Add skip connection, followed by layer normalization
79         return out
80
81 class GPT(nn.Module):
82     def __init__(self, embed_size, num_layers, heads, forward_expansion, dropout, vocab_size):
83         super(GPT, self).__init__()
84         self.embed_size = embed_size
85         self.transformer_blocks = nn.ModuleList(
86             [
87                 TransformerBlock(
88                     embed_size,
89                     heads,
90                     dropout=dropout,
91                     forward_expansion=forward_expansion,
92                 )
93                 for _ in range(num_layers)
94             ]
95         )
96
97         self.word_embedding = nn.Embedding(vocab_size, embed_size)
98         self.position_embedding = nn.Embedding(max_length, embed_size)
99
100     def forward(self, x, mask):
101         N, seq_length = x.shape
102         print(f"Input shape: {x.shape}") # 打印输入形状
103
104         positions = torch.arange(0, seq_length).expand(N, seq_length).to(x.device)
105         out = self.word_embedding(x) + self.position_embedding(positions)
106         print(f"After embedding and position shape: {out.shape}") # 打印嵌入和位置编码后的形状
107
108         for layer in self.transformer_blocks:
109             out = layer(out, out, out, mask)
110             print(f"After transformer block shape: {out.shape}") # 打印经过每个Transformer块后的形状
111
112         return out
113
114 def clean_text_mixed_with_symbols(text):
115     # 保留中文、英文字符、数字和常见的标点符号
116     # 注意：根据需要，你可以在这里添加或删除特定的符号
117     text = re.sub(r'&#91;^\u4e00-\u9fffA-Za-z0-9,.,! ?、;:"\'()《》【】-...]+', ' ', text)
118     return text.strip()
119

```

```

120 def preprocess_text_mixed_with_symbols(text):
121     text = clean_text_mixed_with_symbols(text)
122     tokens = []
123     for token in jieba.cut(text, cut_all=False):
124         token = token.strip()
125         if token:
126             tokens.append(token)
127     return tokens
128
129 def load_and_preprocess_data(file_paths):
130     # 这里简化处理, 具体实现依据你的需求定
131     texts = []
132     for file_path in file_paths:
133         with open(file_path, 'r', encoding='utf-8') as file:
134             text = file.read()
135             # 添加文本清洗和预处理逻辑
136             processed_text = preprocess_text_mixed_with_symbols(text)
137             texts.append(processed_text)
138     return texts
139
140
141 class TextDataset(Dataset):
142     def __init__(self, indexed_texts, vocab_size):
143         self.texts = [torch.tensor(text, dtype=torch.long) for text in indexed_texts]
144         self.vocab_size = vocab_size
145
146     def __len__(self):
147         return len(self.texts)
148
149     def __getitem__(self, idx):
150         return self.texts[idx]
151
152     def collate_fn(self, batch):
153         input_ids = [item[-1] for item in batch]
154         target_ids = [item[1:] for item in batch]
155         input_ids_padded = pad_sequence(input_ids, batch_first=True, padding_value=0)
156         target_ids_padded = pad_sequence(target_ids, batch_first=True, padding_value=0)
157         return input_ids_padded, target_ids_padded
158
159
160 def build_vocab(texts):
161     vocab = set(token for text in texts for token in text)
162     vocab_to_index = {word: i for i, word in enumerate(vocab, start=1)} # 从1开始编号
163     return vocab_to_index
164
165
166 def index_text(text, vocab_to_index):
167     return [vocab_to_index[token] for token in text if token in vocab_to_index]
168
169 def train(model, dataloader, loss_fn, optimizer, device, epochs):
170     model.train()
171     model.to(device)
172
173     for epoch in range(epochs):
174         for batch_idx, (input_ids, target_ids) in enumerate(dataloader):
175             input_ids = input_ids.to(device)
176             target_ids = target_ids.to(device)
177
178             # 前向传播
179             predictions = model(input_ids, mask=None) # 这里简化处理, 没有使用mask
180             predictions = predictions.view(-1, predictions.size(-1))
181             target_ids = target_ids.view(-1)
182
183             # 计算损失
184             loss = loss_fn(predictions, target_ids)

```

```

185
186         # 反向传播和优化
187         optimizer.zero_grad()
188         loss.backward()
189         optimizer.step()
190
191         if batch_idx % 100 == 0:
192             print(f"Epoch {epoch} Batch {batch_idx} Loss {loss.item()}")
193
194
195 # 模型参数
196 vocab_size = 10000 # 假设的词汇表大小
197 embed_size = 256
198 max_length = 100
199 num_layers = 6
200 heads = 8
201 forward_expansion = 4
202 dropout = 0.1
203
204 # 假设你的文本文件路径
205 #file_paths = ['#1';'path/to/your/text1.txt', 'path/to/your/text2.txt']
206 #texts = load_and_preprocess_data(file_paths)
207
208 # 假设文本包含中文、英文和常用符号
209 text = "1977年，三位数学家Rivest、Shamir 和 Adleman 设计了一种算法，可以实现非对称加密。这种算法用他们三
210
211 # 预处理文本
212 texts = preprocess_text_mixed_with_symbols(text)
213 # 输出分词结果
214 print(texts)
215
216 # 假设`texts`是分词后的文本列表
217 vocab_to_index = build_vocab(texts)
218 indexed_texts = ['#1';index_text(text, vocab_to_index) for text in texts]
219
220 vocab_size=len(vocab_to_index) + 1
221 # 现在`texts`应该是索引化后的文本列表
222 dataset = TextDataset(indexed_texts, vocab_size=len(vocab_to_index) + 1) # +1因为从1开始编号
223 dataloader = DataLoader(dataset, batch_size=32, shuffle=True, collate_fn=dataset.collate_fn)
224
225 # 实例化模型
226 model = GPT(
227     embed_size=embed_size,
228     num_layers=num_layers,
229     heads=heads,
230     forward_expansion=forward_expansion,
231     dropout=dropout,
232     vocab_size=vocab_size,
233     max_length=max_length
234 )
235
236 loss_fn = nn.CrossEntropyLoss()
237 optimizer = optim.Adam(model.parameters(), lr=0.0001)
238
239 epochs = 1
240
241 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
242 train(model, dataloader, loss_fn, optimizer, device, epochs)
243
244 # 保存模型参数
245 # model_path = "gpt_simple_model.pth"
246 # torch.save(model.state_dict(), model_path)
247
248 # 如果要保存整个模型（包括模型结构），可以使用以下方式
249 model_path = "gpt_simple_model_full.pth"

```



```
250 torch.save(model, model_path)
251
```

代码提供了一个使用PyTorch实现类似GPT模型的全面示例，这个示例涵盖了多个关键方面，包括自注意力层的定义、变压器块、整体GPT模型、文本数据的预处理（包括混合语言内容的文本清理和使用Jieba进行分词），以及最后的模型训练、自定义数据集和数据加载器的使用。

以下是一些建议和澄清点，以确保代码按预期工作，并遵循最佳实践：

1. **自注意力和变压器块实现**：您的自注意力和变压器块实现看起来很好。它遵循了构建基于变压器模型的标准方法，包括将输入分割成多个头、应用自注意力，然后使用前馈网络。
2. **模型训练循环**：训练循环包括深度学习模型典型训练过程的基本步骤。它通过模型处理输入、计算损失、执行反向传播和更新模型的权重。您还包括了设备兼容性，以便在GPU上运行模型（如果可用），这对于训练效率至关重要。
3. **文本预处理和分词**：您包含了清理文本和分词的功能，这对于NLP任务至关重要。使用Jieba进行分词适用于处理中文文本，您的正则表达式清理混合语言文本涵盖了广泛的字符。
4. **数据处理和数据加载器**：您定义了一个自定义的Dataset类，并使用PyTorch的DataLoader进行批处理和填充。这是处理NLP任务中可变长度序列的好方法。
5. **潜在改进**：
 - **数据预处理中的错误处理**：确保您的文件读取和文本预处理能够优雅地处理错误，尤其是对于可能不存在或有编码问题的文件。
 - **模型中的掩码使用**：您的评论提到了为了简化而没有使用掩码。实际上，特别是对于长度不同的序列，掩码对于通知模型哪些输入部分是填充且不应该被关注至关重要。
 - **词汇表构建**：构建词汇表和索引文本的过程假设所有文本都被分词成一个平面列表。实际上，您可能有多个文档或句子，您可能希望分别处理它们或保持句子边界。
 - **保存模型**：您展示了两种保存模型的方式；仅保存模型参数（state_dict）更节省空间，是大多数用例推荐的方法。保存整个模型虽然方便，但如果需要在不同环境中加载模型，可能会导致问题。

在运行代码之前，请确保调整文件路径，并根据您的具体需求可能扩展预处理和数据集处理。此外，考虑尝试不同的模型超参数（如embed_size、num_layers、heads等）和训练配置，以找到适合您任务的最佳设置。