

基于GPT2搭建医疗问诊机器人

学习目标

- 理解医疗问诊机器人的开发背景.
- 了解企业中聊天机器人的应用场景
- 掌握基于GPT2模型搭建医疗问诊机器人的实现过程

项目背景

- 聊天机器人是一种基于自然语言处理技术的智能对话系统，能够模拟人类的自然语言交流，与用户进行对话和互动。聊天机器人能够理解用户的问题或指令，并给出相应的回答或建议。其目标是提供友好、智能、自然的对话体验。
- 当前，聊天机器人在多个领域得到广泛应用。首先，它们常用于在线客服系统，能够快速、准确地回答用户的常见问题，解决疑问。其次，聊天机器人可以作为个人助手，提供个性化的推荐、建议和日程安排等服务，提升用户体验。此外，聊天机器人还被应用于社交娱乐、语言学习、旅游指南等领域，为用户提供有趣、便捷的对话体验。
- 常见的相关聊天机器人产品：



小冰



小蜜



小度

微软小冰：微软公司开发。它具备自然语言处理、情感分析和对话生成等功能，能够与用户进行智能对话，提供情感支持和娱乐等服务。

阿里云小蜜：阿里云公司推出，提供了丰富的智能对话服务。它具备自然语言处理和对话管理能力，支持多领域的应用场景，如在线客服、智能助手和虚拟导购等。

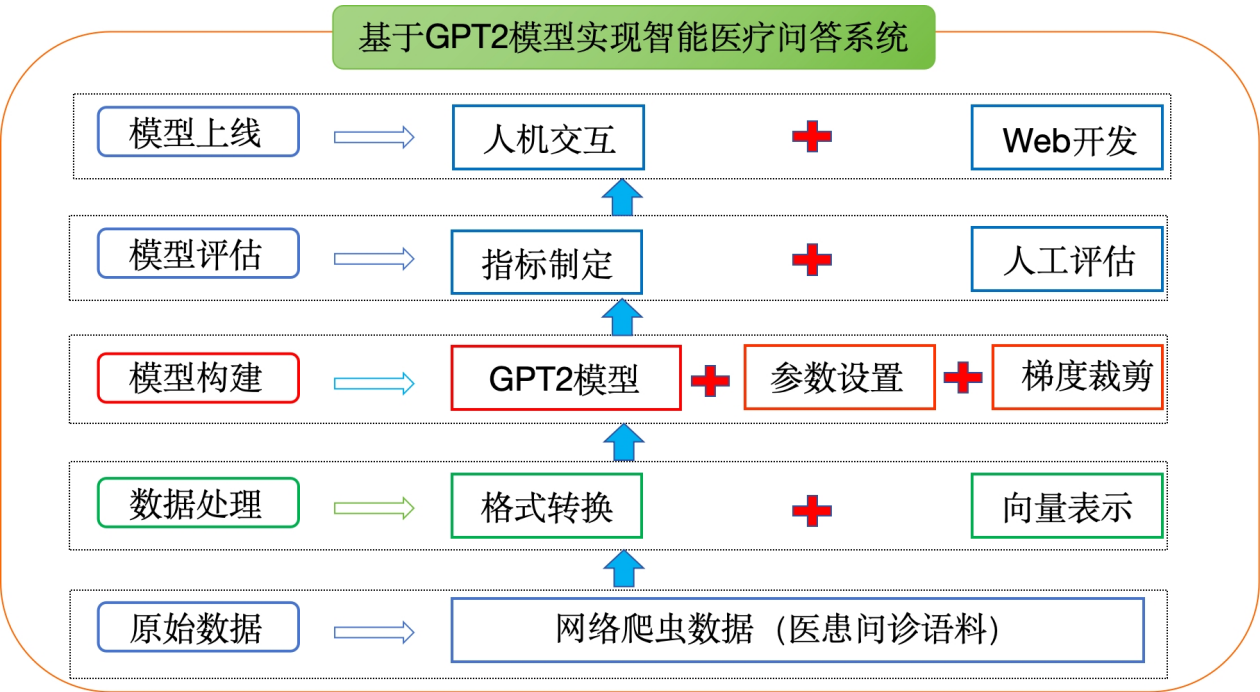
百度智能云小度：百度智能云开发，提供了多领域的智能对话能力。小度机器人可应用于家庭助理、智能音箱和移动应用等场景，通过语音和文本交互与用户进行智能对话，提供信息查询、音乐播放和日程安排等功能。

- 本项目基于医疗领域数据构建了智能医疗问答系统,目的是为用户提供准确、高效、优质的医疗问答服务。

环境准备

- python3.6、
- transformers==4.2.0 、
- pytorch==1.7.0

项目整体结构



1. 数据介绍

- 数据存放位置： /Users/**/PycharmProjects/llm/ptune_chatglm/data
- data文件夹里面包含两个文件： medical_train.txt, medical_valid.txt

1.1 数据展示

- medical_train.txt, medical_valid.txt文件的内容均为对话文本

帕金森叠加综合征的辅助治疗有些什么？

综合治疗；康复训练；生活护理指导；低频重复经颅磁刺激治疗

卵巢癌肉瘤的影像学检查有些什么？

超声漏诊；声像图；MR检查；肿物超声；术前超声；CT检查

低T3综合征的并发症是什么？

心力衰竭；甲状腺结节；糖尿病；感染性休克

重复胚胎停止发育的高危因素有些什么？

黄体功能不足；高龄孕妇

原始train文档中一共包含91487条数据，valid文档中包含1244条数据

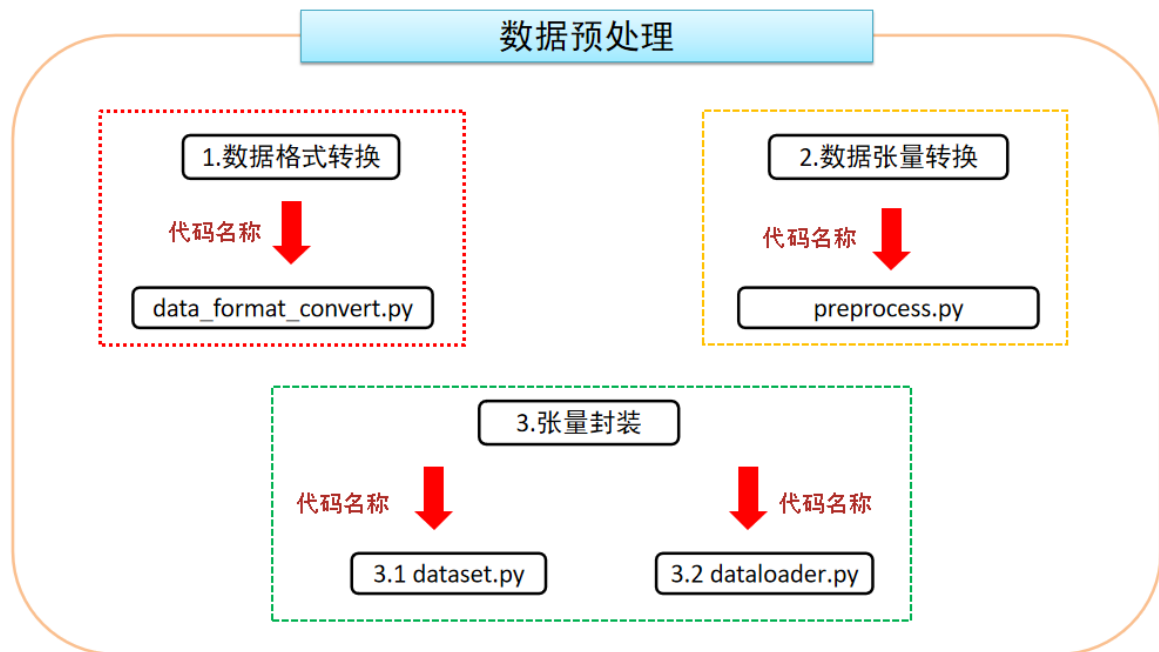
每两行数据为一段对话内容，注意：如果想要实现多轮对话，那么不同对话间实现多条对话语句对即可

2.数据处理

- 目的：将中文文本数据处理成模型能够识别的张量形式。
- 实现过程：
 - 运行preprocess.py，对data/medical_train.txt对话语料进行tokenize，然后进行序列化保存到data/medical_train.txt.pkl。medical_train.pkl中序列化的对象的类型为List[List]，记录对话列表中，每个对话包含的token。

```
1 python preprocess.py
2 train_path data/medical_train.txt --> save_path data/medical_train.pkl
```

- 数据处理基本流程：



2.1 数据张量转换

- 代码路径: `/home/user/ProjectStudy/Gpt2_Chatbot/data_preprocess/preprocess.py`

```
1  # 导入分词器
2  from transformers import BertTokenizerFast
3  # 将数据保存为pkl文件, 方便下次读取
4  import pickle
5  # 读取数据的进度条展示
6  from tqdm import tqdm
7
8
9  def preprocess(train_txt_path, train_pkl_path):
10     """
11     对原始语料进行tokenize, 将每段对话处理成如下形式: "
12     [CLS]utterance1[SEP]utterance2[SEP]utterance3[SEP]"
13     """
14
15     '''初始化tokenizer, 使用BertTokenizerFast.创建一个tokenizer对象'''
16     tokenizer = BertTokenizerFast('./vocab/vocab.txt',
17                                   sep_token="[SEP]",
18                                   pad_token="[PAD]",
19                                   cls_token="[CLS]")
20
21     sep_id = tokenizer.sep_token_id # 获取分隔符[SEP]的token ID
22     cls_id = tokenizer.cls_token_id # 获取起始符[CLS]的token ID
23
24     # 读取训练数据集
25     with open(train_txt_path, 'rb') as f:
26         data = f.read().decode("utf-8") # 以UTF-8编码读取文件内容
```

```

27     # 根据换行符区分不同的对话段落，需要区分Windows和Linux环境下的换行符
28     if "\r\n" in data:
29         train_data = data.split("\r\n\r\n")
30     else:
31         train_data = data.split("\n\n")
32
33     print(len(train_data)) # 打印对话段落数量
34     # 开始进行tokenize
35     # 保存所有的对话数据，每条数据的格式为："[CLS]seq1[SEP]seq2[SEP]seq3[SEP]"
36     dialogue_len = [] # 记录所有对话tokenize之后的长度，用于统计中位数与均值
37     dialogue_list = [] # 记录所有对话
38
39     for index, dialogue in enumerate(tqdm(train_data)):
40         if "\r\n" in data:
41             sequences = dialogue.split("\r\n")
42         else:
43             sequences = dialogue.split("\n")
44
45         input_ids = [cls_id] # 每个dialogue以[CLS]开头
46         for sequence in sequences:
47             # 将每个对话句子进行tokenize，并将结果拼接输入到input_ids列表中
48             input_ids += tokenizer.encode(sequence,
add_special_tokens=False)
49             input_ids.append(sep_id) # 每个seq之后添加[SEP]，表示seqs会话结束
50
51         dialogue_len.append(len(input_ids)) # 将对话的tokenize后的长度添加到
对话长度列表中
52         dialogue_list.append(input_ids) # 将tokenize后的对话添加到对话列表中
53
54     print(f'dialogue_len--->{dialogue_len}') # 打印对话长度列表
55     print(f'dialogue_list--->{dialogue_list}') # 打印
56
57     # 保存pkl文件数据
58     with open(train_pkl_path, "wb") as f:
59         pickle.dump(dialogue_list, f)

```

2.2 数据张量再次封装

2.2.1 封装DataSet对象

- 代码路径：/home/user/ProjectStudy/Gpt2_Chatbot/data_preprocess/dataset.py

```

1  from torch.utils.data import Dataset # 导入Dataset模块，用于定义自定义数据集
2  import torch # 导入torch模块，用于处理张量和构建神经网络
3
4
5  class MyDataset(Dataset):
6      """

```

```

7      自定义数据集类，继承自Dataset类
8      """
9
10     def __init__(self, input_list, max_len):
11         """
12         初始化函数，用于设置数据集的属性
13         :param input_list: 输入列表，包含所有对话的tokenize后的输入序列
14         :param max_len: 最大序列长度，用于对输入进行截断或填充
15         """
16         self.input_list = input_list # 将输入列表赋值给数据集的input_list属性
17         self.max_len = max_len # 将最大序列长度赋值给数据集的max_len属性
18
19     def __len__(self):
20         """
21         获取数据集的长度
22         :return: 数据集的长度
23         """
24         return len(self.input_list) # 返回数据集的长度
25
26     def __getitem__(self, index):
27         """
28         根据给定索引获取数据集中的-一个样本
29         :param index: 样本的索引
30         :return: 样本的输入序列张量
31         """
32         input_ids = self.input_list[index] # 获取给定索引处的输入序列
33         input_ids = input_ids[:self.max_len] # 根据最大序列长度对输入进行截断
或填充
34         input_ids = torch.tensor(input_ids, dtype=torch.long) # 将输入序列
转换为long类型
35         return input_ids # 返回样本的输入序列张量
36

```

2.2.2 封装DataLoader对象

- 代码路径: /home/user/ProjectStudy/Gpt2_Chatbot/data_preprocess/dataloader.py

```

1  # 导入rnn_utils模块，用于处理可变长度序列的填充和排序
2  import torch.nn.utils.rnn as rnn_utils
3
4  # 导入Dataset和DataLoader模块，用于加载和处理数据集
5  from torch.utils.data import Dataset, DataLoader
6
7  import torch # 导入torch模块，用于处理张量和构建神经网络
8  import pickle # 导入pickle模块，用于序列化和反序列化Python对象
9  from dataset import * # 导入自定义的数据集类
10
11  def load_dataset(train_path, valid_path):

```

```

12     # print('进入函数')
13     """
14     加载训练集和验证集
15     :param train_path: 训练数据集路径
16     :return: 训练数据集和验证数据集
17     """
18     with open(train_path, "rb") as f:
19         train_input_list = pickle.load(f) # 从文件中加载输入列表
20
21     with open(valid_path, "rb") as f:
22         valid_input_list = pickle.load(f) # 从文件中加载输入列表
23     # 划分训练集与验证集
24     # print(len(train_input_list)) # 打印输入列表的长度
25     # print(train_input_list[0])
26     #
27     train_dataset = MyDataset(train_input_list, 300) # 创建训练数据集对象
28     val_dataset = MyDataset(valid_input_list, 300) # 创建验证数据集对象
29     return train_dataset, val_dataset # 返回训练数据集和验证数据集
30
31 def collate_fn(batch):
32     """
33     自定义的collate_fn函数，用于将数据集中的样本进行批处理
34     :param batch: 样本列表
35     :return: 经过填充的输入序列张量和标签序列张量
36     """
37     # 对输入序列进行填充，使其长度一致
38     input_ids = rnn_utils.pad_sequence(batch, batch_first=True,
padding_value=0)
39
40     # 对标签序列进行填充，使其长度一致
41     labels = rnn_utils.pad_sequence(batch, batch_first=True,
padding_value=-100)
42
43     return input_ids, labels # 返回经过填充的输入序列张量和标签序列张量
44
45 def get_dataloader(train_path, valid_path):
46     """
47     获取训练数据集和验证数据集的DataLoader对象
48     :param train_path: 训练数据集路径
49     :return: 训练数据集的DataLoader对象和验证数据集的DataLoader对象
50     """
51     # 加载训练数据集和验证数据集
52     train_dataset, val_dataset = load_dataset(train_path, valid_path)
53     print(f'train_dataset-->{len(train_dataset)}')
54     print(f'val_dataset-->{len(val_dataset)}')
55     # 创建训练数据集的DataLoader对象
56     train_dataloader = DataLoader(train_dataset,
57                                     batch_size=4,
58                                     shuffle=True,

```

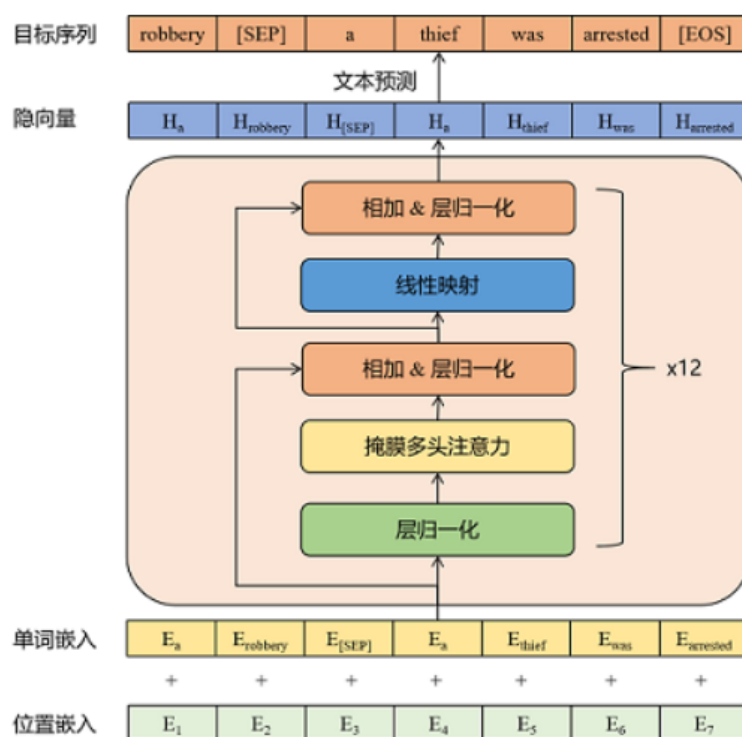
```

59         collate_fn=collate_fn,
60         drop_last=True)
61     # 创建验证数据集的DataLoader对象
62     validate_dataloader = DataLoader(val_dataset,
63                                     batch_size=4,
64                                     shuffle=True,
65                                     collate_fn=collate_fn,
66                                     drop_last=True)
67     # 返回训练数据集的DataLoader对象和验证数据集的DataLoader对象
68     return train_dataloader, validate_dataloader
69
70

```

3. 模型搭建

3.1 模型架构介绍



- 模型架构解析：
 - 输入层：词嵌入层：WordEmbedding + 位置嵌入层：PositionEmbedding
 - 中间层：Transformer的Decoder模块---12层
 - 输出层：LayerNorm层+线性全连接层
- 模型主要参数简介(详见模型的config.json文件):
 - n_embd: 768
 - n_head: 12
 - n_layer: 12
 - n_positions: 1024

- vocab_size: 21128

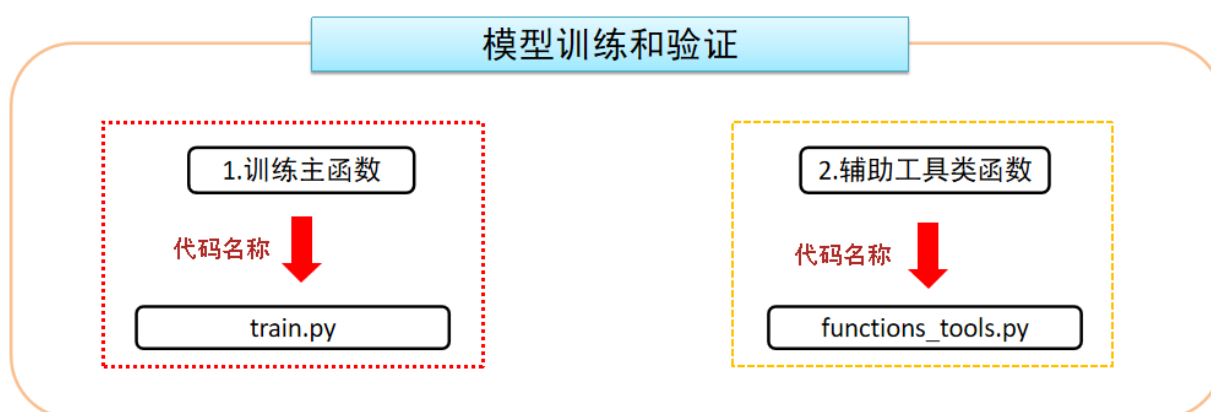
3.2 GPT2模型准备

- 本次项目使用GPT2的预训练模型，因此不需要额外搭建Model类，下面代码是如何直接加载使用GPT2预训练模型
- 代码示例:

```
1 from transformers import GPT2LMHeadModel, GPT2Config
2 # 创建模型
3 if params.pretrained_model:
4     # 加载预训练模型
5     model = GPT2LMHeadModel.from_pretrained(params.pretrained_model)
6 else:
7     # 初始化模型
8     model_config = GPT2Config.from_json_file(params.config_json)
9     model = GPT2LMHeadModel(config=model_config)
```

4. 模型训练和验证

- 主要代码



- 代码路径
 - 训练主函数: /home/user/ProjectStudy/Gpt2_Chatbot/data_preprocess/train.py
 - 辅助工具类
 - 类: /home/user/ProjectStudy/Gpt2_Chatbot/data_preprocess/functions_tools.py

- train.py代码解析

```
1 import torch
2 import os
3 from datetime import datetime
4 import transformers
5 from transformers import GPT2LMHeadModel, GPT2Config
```

```

6  from transformers import BertTokenizerFast
7  from functions_tools import *
8  from parameter_config import *
9  from data_preprocess.dataloader import *
10 from pytorch_tools import EarlyStopping
11
12
13 def train_epoch(model,
14                 train_dataloader,
15                 optimizer, scheduler,
16                 epoch, args):
17     model.train()
18     device = args.device
19     # 对于ignore_index的label token不计算梯度
20     ignore_index = args.ignore_index
21     epoch_start_time = datetime.now()
22     total_loss = 0 # 记录下整个epoch的loss的总和
23
24     # epoch_correct_num:每个epoch中,output预测正确的word的数量
25     # epoch_total_num: 每个epoch中,output预测的word的总数量
26     epoch_correct_num, epoch_total_num = 0, 0
27
28     for batch_idx, (input_ids, labels) in enumerate(train_dataloader):
29         input_ids = input_ids.to(device)
30         labels = labels.to(device)
31         outputs = model.forward(input_ids, labels=labels)
32
33         logits = outputs.logits
34         loss = outputs.loss
35         loss = loss.mean()
36
37
38
39         # 统计该batch的预测token的正确数与总数
40         batch_correct_num, batch_total_num = calculate_acc(logits,
41 labels, ignore_index=ignore_index)
42
43         # 统计该epoch的预测token的正确数与总数
44         epoch_correct_num += batch_correct_num
45         epoch_total_num += batch_total_num
46         # 计算该batch的accuracy
47         batch_acc = batch_correct_num / batch_total_num
48
49         total_loss += loss.item()
50         if args.gradient_accumulation_steps > 1:
51             loss = loss / args.gradient_accumulation_steps
52
53         loss.backward()
54         # 梯度裁剪 # 避免梯度爆炸的方式。

```

```

54         torch.nn.utils.clip_grad_norm_(model.parameters(),
args.max_grad_norm)
55
56         # 进行一定step的梯度累计之后, 更新参数
57         if (batch_idx + 1) % args.gradient_accumulation_steps == 0:
58             # 更新参数
59             optimizer.step()
60             # 更新学习率
61             scheduler.step()
62             # 清空梯度信息
63             optimizer.zero_grad()
64
65             if (batch_idx + 1) % args.loss_step == 0:
66                 print(
67                     "batch {} of epoch {}, loss {}, batch_acc {}, lr
{}".format(
68                         batch_idx + 1, epoch + 1, loss.item() *
args.gradient_accumulation_steps, batch_acc, scheduler.get_lr()))
69
70                 del input_ids, outputs
71
72
73         # 记录当前epoch的平均loss与accuracy
74         epoch_mean_loss = total_loss / len(train_dataloader)
75         epoch_mean_acc = epoch_correct_num / epoch_total_num
76         print(
77             "epoch {}: loss {}, predict_acc {}".format(epoch + 1,
epoch_mean_loss, epoch_mean_acc))
78
79         # save model
80         if epoch % 10 == 0 or epoch == args.epochs:
81             print('saving model for epoch {}'.format(epoch + 1))
82             model_path = os.path.join(args.save_model_path,
'bj_epoch{}'.format(epoch + 1))
83             if not os.path.exists(model_path):
84                 os.mkdir(model_path)
85             model.save_pretrained(model_path)
86             print('epoch {} finished'.format(epoch + 1))
87             epoch_finish_time = datetime.now()
88             print('time for one epoch: {}'.format(epoch_finish_time -
epoch_start_time))
89
90         return epoch_mean_loss
91
92
93 def validate_epoch(model, validate_dataloader, epoch, args):
94     print("start validating")
95     model.eval()
96     device = args.device

```

```

97     ignore_index = args.ignore_index
98     epoch_start_time = datetime.now()
99     total_loss = 0
100     # 捕获cuda out of memory exception
101     with torch.no_grad():
102         for batch_idx, (input_ids, labels) in
enumerate(validate_dataloader):
103             input_ids = input_ids.to(device)
104             labels = labels.to(device)
105             outputs = model.forward(input_ids, labels=labels)
106             logits = outputs.logits
107             loss = outputs.loss
108             loss = loss.mean()
109
110             total_loss += loss.item()
111             del input_ids, outputs
112
113             # 记录当前epoch的平均loss
114             epoch_mean_loss = total_loss / len(validate_dataloader)
115             print(
116                 "validate epoch {}: loss {}".format(epoch+1,
epoch_mean_loss))
117             epoch_finish_time = datetime.now()
118             print('time for validating one epoch:
{}'.format(epoch_finish_time - epoch_start_time))
119             return epoch_mean_loss
120
121
122 def train(model, train_dataloader, validate_dataloader, args):
123
124     # early_stopping = EarlyStopping(patience=0, verbose=True)
125     t_total = len(train_dataloader) // args.gradient_accumulation_steps *
args.epochs
126     optimizer = transformers.AdamW(model.parameters(), lr=args.lr,
eps=args.eps)
127     '''
128     这里对于模型的参数，分别就行权重参数的衰减优化：防止过拟合，以及学习率预热处理优
化：
129     在初始阶段将学习率从较小的值逐步增加到设定的初始值，然后按照设定的学习率调整策略进
行训练。
130     学习率预热的目的是让模型在初始阶段更快地适应数据，避免训练过程中因为学习率过大导致
的梯度爆炸等问题，
131     从而提高模型的训练效果和泛化性能。
132     optimizer: 优化器
133     num_warmup_steps: 初始预热步数
134     num_training_steps: 整个训练过程的总步数
135     '''
136     scheduler = transformers.get_linear_schedule_with_warmup
137     (

```

```

138         optimizer,
139         num_warmup_steps=args.warmup_steps,
140         num_training_steps=t_total
141     )
142
143     print('starting training')
144
145     # 用于记录每个epoch训练和验证的loss
146     train_losses, validate_losses = [], []
147     # 记录验证集的最小loss
148     best_val_loss = 10000
149     # 开始训练
150     for epoch in range(args.epochs):
151         # ===== train ===== #
152         train_loss = train_epoch(
153             model=model, train_dataloader=train_dataloader,
154             optimizer=optimizer, scheduler=scheduler,
155             epoch=epoch, args=args)
156         train_losses.append(train_loss)
157
158         # ===== validate ===== #
159         validate_loss = validate_epoch(
160             model=model, validate_dataloader=validate_dataloader,
161             epoch=epoch, args=args)
162         validate_losses.append(validate_loss)
163
164         # 保存当前困惑度最低的模型，困惑度低，模型的生成效果不一定会越好
165         if validate_loss < best_val_loss:
166             best_val_loss = validate_loss
167             print('saving current best model for epoch {}'.format(epoch +
168 1))
169             model_path = os.path.join(args.save_model_path,
170                                     'min_ppl_model_bj'.format(epoch +
171 1))
172             if not os.path.exists(model_path):
173                 os.mkdir(model_path)
174             model.save_pretrained(model_path)
175
176 def main():
177     # 初始化配置参数
178     params = ParameterConfig()
179
180     # 设置使用哪些显卡进行训练:默认为0
181     os.environ["CUDA_VISIBLE_DEVICES"] = '0'
182
183     # 初始化tokenizer
184     tokenizer = BertTokenizerFast(params.vocab_path,

```

```

185         sep_token="[SEP]",
186         pad_token="[PAD]",
187         cls_token="[CLS]")
188
189
190     # 创建模型的输出目录
191     if not os.path.exists(params.save_model_path):
192         os.mkdir(params.save_model_path)
193
194     # 创建模型
195     if params.pretrained_model:
196         # 加载预训练模型
197         model = GPT2LMHeadModel.from_pretrained(params.pretrained_model)
198     else:
199         # 初始化模型
200         model_config = GPT2Config.from_json_file(params.config_json)
201         model = GPT2LMHeadModel(config=model_config)
202     model = model.to(params.device)
203     assert model.config.vocab_size == tokenizer.vocab_size
204
205
206     # 计算模型参数数量
207     num_parameters = 0
208     parameters = model.parameters()
209     for parameter in parameters:
210         num_parameters += parameter.numel()
211     print(f'模型参数总量---》 {num_parameters}')
212
213     # 加载训练集和验证集
214     # ===== Loading Dataset ===== #
215     train_dataloader, validate_dataloader =
216     get_dataloader(params.train_path)
217     train(model, train_dataloader, validate_dataloader, params)
218
219 if __name__ == '__main__':
220     main()

```

- functions_tools.py代码解析

```

1  import torch
2  import torch.nn.functional as F
3
4  def calculate_acc(logits, labels, ignore_index=-100):
5      logits = logits[:, :-1, :].contiguous().view(-1, logits.size(-1))
6      labels = labels[:, 1:].contiguous().view(-1)
7      _, logits = logits.max(dim=-1) # 对于每条数据, 返回最大的index
8      '''

```

```

9      在 PyTorch 中, labels.ne(ignore_index) 表示将标签张量 labels 中的值不等于
      ignore_index 的位置标记为 True, 等于 ignore_index 的位置标记为 False。
10     这个操作通常被用于计算交叉熵损失, 以过滤掉 ignore_index 对损失的贡献
11     '''
12     # 进行非运算, 返回一个tensor, 若labels的第i个位置为pad_id, 则置为0, 否则为1
13     non_pad_mask = labels.ne(ignore_index)
14     '''
15     在 PyTorch 中,
16     logit.eq(labels) 表示将模型的预测输出值 logit 中等于标签张量 labels 的位置标记
      为 True, 不等于标签张量 labels 的位置标记为 False。这个操作通常被用于计算交叉熵损失,
      以标记出预测输出值和标签值相等的位置。
17     masked_select(non_pad_mask) 表示将张量中非填充标记的位置选出来。这个操作通常被
      用于计算损失时, 过滤掉填充标记对损失的影响。
18     '''
19     n_correct = logit.eq(labels).masked_select(non_pad_mask).sum().item()
20     n_word = non_pad_mask.sum().item()
21     return n_correct, n_word

```

5. 模型预测（人机交互）

- 运行interact.py, 使用训练好的模型, 进行人机交互, 输入Ctrl+Z结束对话之后, 聊天记录将保存到sample目录下的sample.txt文件中。

```

1  import os
2  from datetime import datetime
3  from transformers import GPT2LMHeadModel
4  from transformers import BertTokenizerFast
5  import torch.nn.functional as F
6  from parameter_config import *
7
8  PAD = '[PAD]'
9  pad_id = 0
10
11
12  def top_k_top_p_filtering(logits, top_k=0, top_p=0.0, filter_value=-
      float('Inf')):
13      """使用top-k和/或nucleus (top-p) 筛选来过滤logits的分布
14      参数:
15          logits: logits的分布, 形状为 (词汇大小)
16          top_k > 0: 保留概率最高的top k个标记 (top-k筛选)。
17          top_p > 0.0: 保留累积概率大于等于top_p的top标记 (nucleus筛选)。
18      """
19      assert logits.dim() == 1 # batch size 1 for now - could be updated
      for more but the code would be less clear
20      top_k = min(top_k, logits.size(-1)) # Safety check
21      print(f'top_k---->{top_k}')
22      if top_k > 0:

```

```

23         # Remove all tokens with a probability less than the last token
of the top-k
24         # torch.topk()返回最后一维最大的top_k个元素, 返回值为二维
(values,indices)
25         # ...表示其他维度由计算机自行推断
26         print(f'torch.topk(logits, top_k)-->{torch.topk(logits, top_k)}')
27         indices_to_remove = logits < torch.topk(logits, top_k)[0][...,
-1, None]
28         logits[indices_to_remove] = filter_value # 对于topk之外的其他元素的
logits值设为负无穷
29
30         if top_p > 0.0:
31             sorted_logits, sorted_indices = torch.sort(logits,
descending=True) # 对logits进行递减排序
32             print(f'sorted_logits-->{sorted_logits}')
33             print(f'sorted_indices-->{sorted_indices}')
34             cumulative_probs = torch.cumsum(F.softmax(sorted_logits, dim=-1),
dim=-1)
35
36             # Remove tokens with cumulative probability above the threshold
37             sorted_indices_to_remove = cumulative_probs > top_p
38             # Shift the indices to the right to keep also the first token
above the threshold
39             sorted_indices_to_remove[..., 1:] = sorted_indices_to_remove[...,
:-1].clone()
40             sorted_indices_to_remove[..., 0] = 0
41
42             indices_to_remove = sorted_indices[sorted_indices_to_remove]
43             logits[indices_to_remove] = filter_value
44         return logits
45
46
47 def main():
48     pconf = ParameterConfig()
49     # 当用户使用GPU,并且GPU可用时
50     device = 'cuda' if torch.cuda.is_available() else 'cpu'
51     print('using device:{}'.format(device))
52     os.environ["CUDA_VISIBLE_DEVICES"] = '0'
53     tokenizer = BertTokenizerFast(vocab_file=pconf.vocab_path,
54                                 sep_token="[SEP]",
55                                 pad_token="[PAD]",
56                                 cls_token="[CLS]")
57     model = GPT2LMHeadModel.from_pretrained('./save_model/epoch25')
58     model = model.to(device)
59     model.eval()
60     # 保存聊天记录的文件路径
61     if pconf.save_samples_path:
62         if not os.path.exists(pconf.save_samples_path):
63             os.makedirs(pconf.save_samples_path)

```



```

64     samples_file = open(pconf.save_samples_path + '/samples.txt',
'a', encoding='utf8')
65     samples_file.write("聊天记录{}:\n".format(datetime.now()))
66     # 存储聊天记录, 每个utterance以token的id的形式进行存储
67     history = []
68     print('开始和chatbot聊天, 输入CTRL + z以退出')
69
70     while True:
71         try:
72             text = input("user:")
73             # text = "你好"
74             if pconf.save_samples_path:
75                 samples_file.write("user:{}\n".format(text))
76             text_ids = tokenizer.encode(text, add_special_tokens=False)
77             print(f'text_ids-->{text_ids}')
78             print('*' * 80)
79             history.append(text_ids)
80             input_ids = [tokenizer.cls_token_id] # 每个input以[CLS]为开头
81             print(f'history---.{history}')
82             print(f'input_ids---.{input_ids}')
83             print('*' * 80)
84             print(f'history[-pconf.max_history_len:]-->{history[-
pconf.max_history_len:]}')
85             for history_id, history_utr in enumerate(history[-
pconf.max_history_len:]):
86                 input_ids.extend(history_utr)
87                 print(input_ids)
88                 input_ids.append(tokenizer.sep_token_id)
89                 print(input_ids)
90             print('*'*80)
91             print(f'new_inut--->{input_ids}')
92             input_ids = torch.tensor(input_ids).long().to(device)
93             input_ids = input_ids.unsqueeze(0)
94             print(f'las--inputs_ids{input_ids}')
95             response = [] # 根据context, 生成的response
96             # 最多生成max_len个token
97             for _ in range(pconf.max_len):
98                 outputs = model(input_ids=input_ids)
99                 logits = outputs.logits
100                 print(f'logits--->{logits}')
101                 print(f'logits--->{logits.shape}')
102                 print('*'*80)
103                 next_token_logits = logits[0, -1, :]
104                 # 对于已生成的结果generated中的每个token添加一个重复惩罚项, 降低
其生成概率
105                 print(f'next_token_logits-->{next_token_logits}')
106                 for id in set(response):
107                     print(f'id--->{id}')
108                     next_token_logits[id] /= pconf.repetition_penalty

```

```

109         # 对于[UNK]的概率设为无穷小, 也就是说模型的预测结果不可能是[UNK]这
        个token
110
111         next_token_logits[tokenizer.convert_tokens_to_ids('[UNK]')] = -
        float('Inf')
112         filtered_logits =
        top_k_top_p_filtering(next_token_logits, top_k=pconf.topk,
        top_p=pconf.topp)
113         print(f'filtered_logits-->{filtered_logits}')
114         # torch.multinomial表示从候选集中无放回地进行抽取num_samples
        个元素, 权重越高, 抽到的几率越高, 返回元素的下标
115         next_token = torch.multinomial(F.softmax(filtered_logits,
        dim=-1), num_samples=1)
116         print(f'next_token-->{next_token}')
117         if next_token == tokenizer.sep_token_id: # 遇到[SEP]则表
        明response生成结束
118             break
119             response.append(next_token.item())
120             input_ids = torch.cat((input_ids,
        next_token.unsqueeze(0)), dim=1)
121             # his_text =
        tokenizer.convert_ids_to_tokens(curr_input_tensor.tolist())
122             # print("his_text:{}".format(his_text))
123             print(f'response-->{response}')
124             history.append(response)
125             text = tokenizer.convert_ids_to_tokens(response)
126             print("chatbot:" + "".join(text))
127             if pconf.save_samples_path:
128                 samples_file.write("chatbot:{}\n".format("".join(text)))
129             except KeyboardInterrupt:
130                 if pconf.save_samples_path:
131                     samples_file.close()
132                 break
133
134 if __name__ == '__main__':
135     main()

```