

基于BERT+PET方式数据预处理介绍

学习目标

- 了解本项目数据类型和表现格式
- 掌握数据处理的工具函数代码实现

BERT+PET方式数据预处理

- 本项目中对数据部分的预处理步骤如下:

1. 查看项目数据集
2. 编写Config类项目文件配置代码
3. 编写数据处理相关代码

1 查看项目数据集

- 数据存放位置: /Users/***/PycharmProjects/llm/prompt_tasks/PET/data
- data文件夹里面包含4个txt文档, 分别为: train.txt、dev.txt、prompt.txt、verbalizer.txt

1.1 train.txt

- train.txt为训练数据集, 其部分数据展示如下:

1	水果	脆脆的, 甜味可以, 可能时间有点长了, 水分不是很足。
2	平板	华为机器肯定不错, 但第一次碰上京东最糟糕的服务, 以后没想到京东购物了。
3	书籍	为什么不认真的检查一下, 发这么一本脏脏的书给顾客呢!
4	衣服	手感不错, 用料也很好, 不知道水洗后怎样, 相信大品牌, 质量过关, 五星好评!!!
5	水果	苹果有点小, 不过好吃, 还有几个烂的。估计是故意的放的。差评。
6	衣服	掉色掉的厉害, 洗一次就花了

train.txt一共包含63条样本数据, 每一行用\t 分开, 前半部分为标签(label), 后半部分为原始输入(用户评论)。

如果想使用自定义数据训练, 只需要仿照上述示例数据构建数据集即可。

1.2 dev.txt

- dev.txt为验证数据集, 其部分数据展示如下:

1	书籍	"一点都不好笑,很失望,内容也不是很实用"
2	衣服	完全是一条旧裤子。
3	手机	相机质量不错,如果阳光充足,可以和数码相机媲美。界面比较人性化,容易使用。软件安装简便
4	书籍	明明说有货,结果送货又没有了。并且也不告诉我,怎么评啊
5	洗浴	非常不满意,晚上洗的头发,第二天头痒痒的不行了,还都是头皮屑。
6	水果	这个苹果感觉是长熟的苹果,没有打蜡,不错,又甜又脆

dev.txt一共包含590条样本数据，每一行用\t分开，前半部分为标签(label)，后半部分为原始输入(用户评论)。

如果想使用自定义数据训练，只需要仿照上述示例数据构建数据集即可。

1.3 prompt.txt

- prompt.txt为人工设定提示模版，其数据展示如下：

1 | 这是一条{MASK}评论：{textA}。

其中，用大括号括起来的部分为「自定义参数」，可以自定义设置大括号内的值。

示例中{MASK}代表[MASK] token的位置，{textA}代表评论数据的位置。

你可以改为自己想要的模板，例如想新增一个{textB}参数：

1 | {textA}和{textB}是{MASK}同的意思。

1.4 verbalizer.txt

- verbalizer.txt主要用于定义「真实标签」到「标签预测词」之间的映射。在有些情况下，将「真实标签」作为[MASK]去预测可能不具备很好的语义通顺性，因此，我们会对「真实标签」做一定的映射。
- 例如：

1 | "中国爆冷2-1战胜韩国"是一则[MASK][MASK]新闻。 体育

- 这句话中的标签为「体育」，但如果我们将标签设置为「足球」会更容易预测。
- 因此，我们可以对「体育」这个label构建许多个子标签，在推理时，只要预测到子标签最终推理出真实标签即可，如下：

1 | 体育 -> 足球,篮球,网球,棒球,乒乓,体育

- 项目中标签词映射数据展示如下：

1	电脑	电脑
2	水果	水果
3	平板	平板
4	衣服	衣服
5	酒店	酒店
6	洗浴	洗浴
7	书籍	书籍
8	蒙牛	蒙牛
9	手机	手机
10	电器	电器

verbalizer.txt 一共包含10个类别，上述数据中，我们使用了1对1的verbalizer, 如果想定义一对多的映射，只需要在后面用","分割即可， eg:

```
1 | 水果 苹果,香蕉,橘子
```

若想使用自定义数据训练，只需要仿照示例数据构建数据集

2 编写Config类项目文件配置代码

- 代码路径: /Users/***/PycharmProjects/lm/prompt_tasks/PET/pet_config.py
- config文件目的: 配置项目常用变量，一般这些变量属于不经常改变的，比如：训练文件路径、模型训练次数、模型超参数等等

具体代码实现：

```
1 # coding:utf-8
2 import torch
3 import sys
4 print(sys.path)
5
6 class ProjectConfig(object):
7     def __init__(self):
8         # 是否使用GPU
9         self.device = 'cuda:0' if torch.cuda.is_available() else 'cpu'
10        # 预训练模型bert路径
11        self.pre_model = '/home/prompt_project/bert-base-chinese'
12        self.train_path = '/home/prompt_project/PET/data/train.txt'
13        self.dev_path = '/home/prompt_project/PET/data/dev.txt'
14        self.prompt_file = '/home/prompt_project/PET/data/prompt.txt'
15        self.verbalizer = '/home/prompt_project/PET/data/verbalizer.txt'
16        self.max_seq_len = 512
17        self.batch_size = 8
18        self.learning_rate = 5e-5
19        # 权重衰减参数（正则化，抑制模型过拟合）
20        self.weight_decay = 0
21        # 预热学习率(用来定义预热的步数)
22        self.warmup_ratio = 0.06
```

```

23     self.max_label_len = 2
24     self.epochs = 50
25     self.logging_steps = 10
26     self.valid_steps = 20
27     self.save_dir = '/home/prompt_project/PET/checkpoints'
28
29
30 if __name__ == '__main__':
31     pc = ProjectConfig()
32     print(pc.prompt_file)
33     print(pc.pre_model)

```

3 编写数据处理相关代码

- 代码路径: /Users/***/PycharmProjects/llm/prompt_tasks/PET/data_handle.
- data_handle文件夹中一共包含三个py脚本: template.py、data_preprocess.py、data_loader.py

3.1 template.py

- 目的: 构建固定模版类, text2id的转换
- 导入必备工具包

```

1 # -*- coding:utf-8 -*-
2 from rich import print # 终端层次显示
3 from transformers import AutoTokenizer
4 import numpy as np
5 from pet_config import *

```

- 定义HardTemplate类

```

1 class HardTemplate(object):
2     """
3         硬模板，人工定义句子和[MASK]之间的位置关系。
4     """
5
6     def __init__(self, prompt: str):
7         """
8             Args:
9                 prompt (str): prompt格式定义字符串，e.g. -> "这是一条{MASK}评论:
10 {textA}。"
11
12             self.prompt = prompt
13             self.inputs_list = [] # 根据文字prompt拆分为
14             # 各part的列表
15             self.custom_tokens = set(['MASK']) # 从prompt中解析出的自
16             # 定义token集合
17             self.prompt_analysis() # 解析prompt模板

```

```
15
16     def prompt_analysis(self):
17         """
18             将prompt文字模板拆解为可映射的数据结构。
19
20             Examples:
21                 prompt -> "这是一条{MASK}评论: {textA}。"
22                 inputs_list -> ['这', '是', '一', '条', 'MASK', '评', '论', ': ', 'textA', '.']
23                 custom_tokens -> {'textA', 'MASK'}
24             """
25
26             idx = 0
27             while idx < len(self.prompt):
28                 str_part = ''
29                 if self.prompt[idx] not in ['{', '}']:
30                     self.inputs_list.append(self.prompt[idx])
31                 if self.prompt[idx] == '{':                      # 进入自定义字段
32                     idx += 1
33                     while self.prompt[idx] != '}':
34                         str_part += self.prompt[idx]                  # 拼接该自定义
35
36             idx += 1
37             elif self.prompt[idx] == '}':
38                 raise ValueError("Unmatched bracket '{}', check your
39 prompt.")
38             if str_part:
39                 self.inputs_list.append(str_part)
40                 # 将所有自定义字段存储, 后续会检测输入信息是否完整
41                 self.custom_tokens.add(str_part)
42             idx += 1
43
44     def __call__(self,
45                 inputs_dict: dict,
46                 tokenizer,
47                 mask_length,
48                 max_seq_len=512):
49             """
50                 输入一个样本, 转换为符合模板的格式。
51
52             Args:
53                 inputs_dict (dict): prompt中的参数字典, e.g. -> {
54                     "textA": "这
55                     个手机也太卡了",
56                     "[MASK]"
57
58             }
59                 tokenizer: 用于encoding文本
60                 mask_length (int): MASK token 的长度
```

```
59     Returns:
60         dict -> {
61             'text': '[CLS]这是一条[MASK]评论：这个手机也太卡了。[SEP]',
62             'input_ids': [1, 47, 10, 7, 304, 3, 480, 279, 74, 47,
63                         27, 247, 98, 105, 512, 777, 15, 12043, 2],
64             'token_type_ids': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
65                                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
66             'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
67                                 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],
68             'mask_position': [0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
69                                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
70         }
71     """
72     # 定义输出格式
73     outputs = {
74         'text': '',
75         'input_ids': [],
76         'token_type_ids': [],
77         'attention_mask': [],
78         'mask_position': []
79     }
80
81     str_formated = ''
82     for value in self.inputs_list:
83         if value in self.custom_tokens:
84             if value == 'MASK':
85                 str_formated += inputs_dict[value] * mask_length
86             else:
87                 str_formated += inputs_dict[value]
88         else:
89             str_formated += value
90     # print(f'str_formated-->{str_formated}')
91     encoded = tokenizer(text=str_formated,
92                         truncation=True,
93                         max_length=max_seq_len,
94                         padding='max_length')
95     # print(f'encoded--->{encoded}')
96     outputs['input_ids'] = encoded['input_ids']
97     outputs['token_type_ids'] = encoded['token_type_ids']
98     outputs['attention_mask'] = encoded['attention_mask']
99     token_list =
100     tokenizer.convert_ids_to_tokens(encoded['input_ids'])
101     outputs['text'] = ''.join(token_list)
102     mask_token_id = tokenizer.convert_tokens_to_ids(['[MASK]'])[0]
103     condition = np.array(outputs['input_ids']) == mask_token_id
104     mask_position = np.where(condition)[0].tolist()
105     outputs['mask_position'] = mask_position
106     return outputs
```

```

107
108 if __name__ == '__main__':
109     pc = ProjectConfig()
110     tokenizer = AutoTokenizer.from_pretrained(pc.pre_model)
111     hard_template = HardTemplate(prompt='这是一条{MASK}评论: {textA}')
112     print(hard_template.inputs_list)
113     print(hard_template.custom_tokens)
114     tep = hard_template(
115         inputs_dict={'textA': '包装不错，苹果挺甜的，个头也大。',
116                     'MASK': '[MASK]'},
117         tokenizer=tokenizer,
118         max_seq_len=30,
119         mask_length=2)
120     print(tep)
121
122     print(tokenizer.convert_ids_to_tokens([3819, 3352]))
123     print(tokenizer.convert_tokens_to_ids(['水', '果']))

```

3.2 data_preprocess.py

- 目的: 将样本数据转换为模型接受的输入数据
- 导入必备的工具包

```

1 from template import *
2 from rich import print
3 from datasets import load_dataset
4 # partial: 把一个函数的某些参数给固定住（也就是设置默认值），返回一个新的函数，调用这个新
      函数会更简单
5 from functools import partial
6 from pet_config import *

```

- 定义数据转换方法convert_example()

```

1 def convert_example(
2     examples: dict,
3     tokenizer,
4     max_seq_len: int,
5     max_label_len: int,
6     hard_template: HardTemplate,
7     train_mode=True,
8     return_tensor=False) -> dict:
9     """
10     将样本数据转换为模型接收的输入数据。
11
12     Args:
13         examples (dict): 训练数据样本, e.g. -> {

```

```
14             "text": [
15                 '手机 这个手机
16                 也太卡了。',
17                 '体育 世界杯为
18                 何迟迟不见宣传',
19                 ...
20             ]
21         }
22     max_seq_len (int): 句子的最大长度, 若没有达到最大长度, 则padding为最大长
度
23     max_label_len (int): 最大label长度, 若没有达到最大长度, 则padding为最大
长度
24     hard_template (HardTemplate): 模板类。
25     train_mode (bool): 训练阶段 or 推理阶段。
26     return_tensor (bool): 是否返回tensor类型, 如不是, 则返回numpy类型。
27
28     Returns:
29         dict (str: np.array) -> tokenized_output = {
30             'input_ids': [[1, 47, 10, 7, 304, 3, 3, 3, 3,
31             47, 27,
32             247, 98, 105, 512, 777, 15,
33             12043, 2], ...],
34             'token_type_ids': [[0, 0, 0, 0, 0, 0, 0, 0, 0,
35             0, 0, 0,
36             0, 0, 0, 0, 0, 0, 0, 0, 0],
37             ...],
38             'attention_mask': [[1, 1, 1, 1, 1, 1, 1, 1, 1,
39             1, 1, 1,
40             1, 1, 1, 1, 1, 1, 1, 1],
41             ...],
42             'mask_positions': [[5, 6, 7, 8], ...],
43             'mask_labels': [[2372, 3442, 0, 0],
44             [2643, 4434, 2334, 0], ...]
45         }
46
47         """
48         tokenized_output = {
49             'input_ids': [],
50             'token_type_ids': [],
51             'attention_mask': [],
52             'mask_positions': [],
53             'mask_labels': []
54         }
55
56         for i, example in enumerate(examples['text']):
57             if train_mode:
58                 label, content = example.strip().split('\t')
59             else:
60                 content = example.strip()
```

```

53     inputs_dict = {
54         'textA': content,
55         'MASK': '[MASK]'
56     }
57     encoded_inputs = hard_template(
58         inputs_dict=inputs_dict,
59         tokenizer=tokenizer,
60         max_seq_len=max_seq_len,
61         mask_length=max_label_len)
62     tokenized_output['input_ids'].append(encoded_inputs["input_ids"])
63
64     tokenized_output['token_type_ids'].append(encoded_inputs["token_type_ids"])
65
66     tokenized_output['attention_mask'].append(encoded_inputs["attention_mask"])
67
68     tokenized_output['mask_positions'].append(encoded_inputs["mask_position"])
69
70     if train_mode:
71         label_encoded = tokenizer(text=[label]) # 将label补到最大长度
72         # print(f'label_encoded-->{label_encoded}')
73         label_encoded = label_encoded['input_ids'][0][1:-1]
74         label_encoded = label_encoded[:max_label_len]
75         add_pad = [tokenizer.pad_token_id] * (max_label_len -
76         len(label_encoded))
77         label_encoded = label_encoded + add_pad
78         tokenized_output['mask_labels'].append(label_encoded)
79
80     for k, v in tokenized_output.items():
81         if return_tensor:
82             tokenized_output[k] = torch.LongTensor(v)
83         else:
84             tokenized_output[k] = np.array(v)
85
86     return tokenized_output
87
88
89
90
91
92
93
94

```

```
95     convert_func = partial(convert_example,
96                             tokenizer=tokenizer,
97                             hard_template=hard_template,
98                             max_seq_len=30,
99                             max_label_len=2)
100    dataset = train_dataset.map(convert_func, batched=True)
101    for value in dataset['train']:
102        print(value)
103        print(len(value['input_ids']))
104        break
105
```

3.3 data_loader.py

- 目的：定义数据加载器
- 导入必备的工具包

```
1 # coding:utf-8
2 from torch.utils.data import DataLoader
3 from transformers import default_data_collator
4 from data_preprocess import *
5 from pet_config import *
6
7 pc = ProjectConfig() # 实例化项目配置文件
8 tokenizer = AutoTokenizer.from_pretrained(pc.pre_model)
```

- 定义获取数据加载器的方法get_data()

```
1 def get_data():
2     # prompt定义
3     prompt = open(pc.prompt_file, 'r', encoding='utf8').readlines()
4     [0].strip()
5     hard_template = HardTemplate(prompt=prompt) # 模板转换器定义
6     dataset = load_dataset('text', data_files={'train': pc.train_path,
7                                                 'dev': pc.dev_path})
8     # print(dataset)
9     # print(f'Prompt is -> {prompt}')
10    new_func = partial(convert_example,
11                        tokenizer=tokenizer,
12                        hard_template=hard_template,
13                        max_seq_len=pc.max_seq_len,
14                        max_label_len=pc.max_label_len)
15
16    dataset = dataset.map(new_func, batched=True)
17
18    train_dataset = dataset["train"]
```

```
18 dev_dataset = dataset[ "dev" ]
19 # print('train_dataset', train_dataset[:2])
20 # print('*'*80)
21 train_dataloader = DataLoader(train_dataset,
22                             shuffle=True,
23                             collate_fn=default_data_collator,
24                             batch_size=pc.batch_size)
25 dev_dataloader = DataLoader(dev_dataset,
26                             collate_fn=default_data_collator,
27                             batch_size=pc.batch_size)
28 return train_dataloader, dev_dataloader
29
30
31 if __name__ == '__main__':
32     train_dataloader, dev_dataloader = get_data()
33     print(len(train_dataloader))
34     print(len(dev_dataloader))
35     for i, value in enumerate(train_dataloader):
36         print(i)
37         print(value)
38         print(value[ 'input_ids' ].dtype)
39         break
40
```