

# 基于 Bert 实现临床试验筛选标准短文本分类 (CHIP-CTC)

马英梓

2022 年 5 月 29 日

## 1 前言

本次实验基于中文医疗信息处理挑战榜 CBLUE(Chinese Biomedical Language Understanding Evaluation), 根据数据集 CHIP-CTC(CHiP - Clinical Trial Criterion dataset) 完成临床试验筛选标准短文本分类任务, 数据量为: 训练集数据 22962 条, 验证集数据 7682 条, 测试集数据 10000 条。。在 baseline 模型的选择中, 我选择使用 BERT-www-ext-base 模型 [1], 这是一个由哈工大讯飞联合实验室推出的基于全词掩码 (Whole Word Masking) 技术的中文预训练模型, 配置为 12-layer, 768-hidden, 12-heads, 110M parameters。实验的目的旨在验证 NLP 下游任务训练中的一些 tricks 在该任务中是否 work。

## 2 任务

### 2.1 任务介绍

给定事先定义好的 44 种筛选标准语义类别和一系列中文临床试验筛选标准的描述句子, 参赛者需返回每一条筛选标准的具体类别。

ID	输入 (筛选标准)	输出 (类别)
迁移	Text(input)	Label(output)
探究	年龄 >80 岁	Age
融合	近期颅内或椎管内手术史	Therapy or Surgery
展示	血糖 <2.7mmol/L	Laboratory Examinations

表 1: 标注数据示例

查看数据集 label 分布, 根据图 1 可以明显看出数据集中存在样本类别分布不均匀的现象, 因此想要取得一个不错的结果需要首先解决样本类别分布不均匀的问题。因此, 我打算对不同 label 的样本做不同数量的数据增强, 尽可能平衡样本数量。

### 2.2 评测指标

本任务的评价指标使用宏观 F1 值 (Macro-F1, 或称 Average-F1)。该指标先计算每一类的准确率和召回率, 然后对所有类别的准确率和召回率求均值, 最后后应用 F1 计算公式。因为对各类别的 Precision 和 Recall 求了平均, 所以并没有考虑到数据数量的问题。在这种情况下, Precision 和 Recall 较高的类别对 F1 的影响会较大。

计算公式简化为:

$$Precision_{macro} = \frac{\sum_{i=1}^n Precision_i}{n}$$

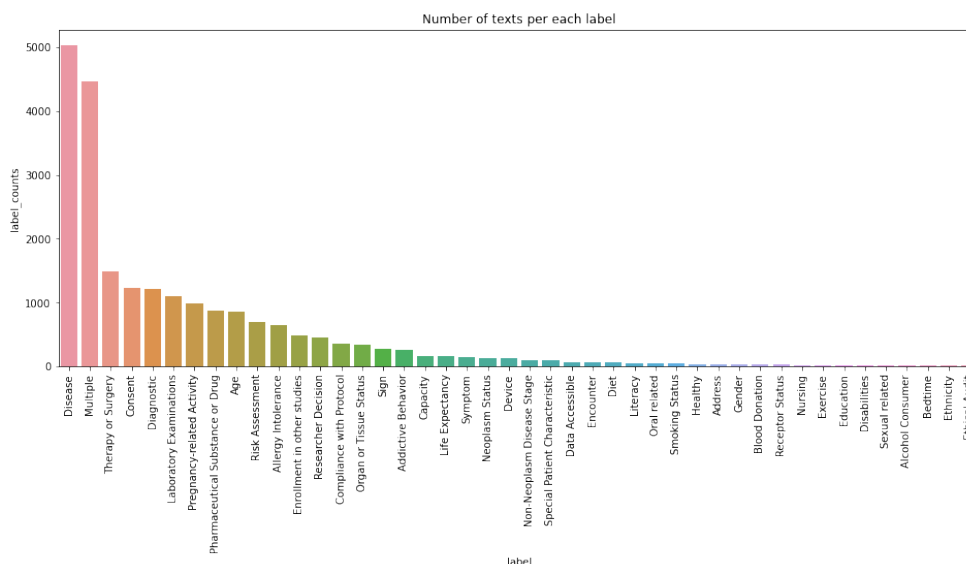


图 1: 数据集中 label 的数量分布

$$Recall_{macro} = \frac{\sum_{i=1}^n Recall_i}{n}$$

$$Macro - F1 = 2 \times \frac{Precision_{macro} * Recall_{macro}}{Precision_{macro} + Recall_{macro}}$$

## 2.3 Tricks

### 2.3.1 Data Preprocessing

检索数据时发现，每条数据前都有类似 (1)、a) 这种标号，而这些标号对于文本的分类是没有影响的，因此利用正则表达式去掉。

```
1 a = re.findall('[\u4e00-\u9fa5A-Za-z][\S\s]+', s, re.S) #只要字符串中的中文，字母，数字
2 a = "".join(a)
3
4 output:
5 (4) 初诊及复发患者且6个月内未经放，化疗诊治。
6 初诊及复发患者且6个月内未经放，化疗诊治。
```

### 2.3.2 Chinese Data Augmentation

这里我使用了 nlpca 库来对训练文本做数据增强。我从中选择了 3 种方法对文本进行处理。

- 随机同义词替换：

```
1 from nlpca import Similarword
2 smw = Similarword(create_num=3, change_rate=0.3)
3 rs1 = smw.replace(a)
4
5 output:
6 已知过敏体质的患者，如对人血白蛋白过敏者,对手术过所需任何药品试剂过敏者;
7 随机同义词替换>>>>>>
8 已知过敏体质的病家，如对人血白蛋白过敏者,对手术过所需另药品试剂过敏者;
9 已知过敏体质的患者，如对人血白蛋白过敏者,对手术过所需另药味试剂过敏者;
```

- 随机字删除，同时也能去除尾部的一些标点噪声：

```
1 from nlpcda import RandomDeleteChar
2 smw = RandomDeleteChar(create_num=3, change_rate=0.3)
3 rs1 = smw.replace(a)

5 output :
有服用镇静剂、抗失眠药物及抗抑郁药的病史；
7 随机字删除>>>>>
服用镇静剂、抗失眠药物抗抑郁药的病史，
9 有服用镇静剂失眠药物及抗抑郁药的病史
```

- 翻译互转实现的增强，利用百度翻译来回翻译清洗数据：

```
1 from nlpcda import baidu_translate
2 en_s = baidu_translate(content=a, appid='20220528001232402', secretKey='6yH2KBcds7tNR1RbynpJ', t_from='zh', t_to='en')
3 zh_s = baidu_translate(content=en_s, appid='20220528001232402', secretKey='6yH2KBcds7tNR1RbynpJ', t_from='en', t_to='zh')

5 output :
已知过敏体质的患者，如对人血白蛋白过敏者,对手术过所需任何药品试剂过敏者；
7 翻译互转>>>>>
Patients with known allergic constitution, such as those who are allergic to human albumin and those
9 who are allergic to any drugs and reagents required for surgery;
具有已知过敏体质的患者，如对人类白蛋白过敏的患者和对手术所需的任何药物和试剂过敏的患者；
```

在使用了3种数据增强后，我尽可能调整了样本的label分布，调整后的结果如图2所示。

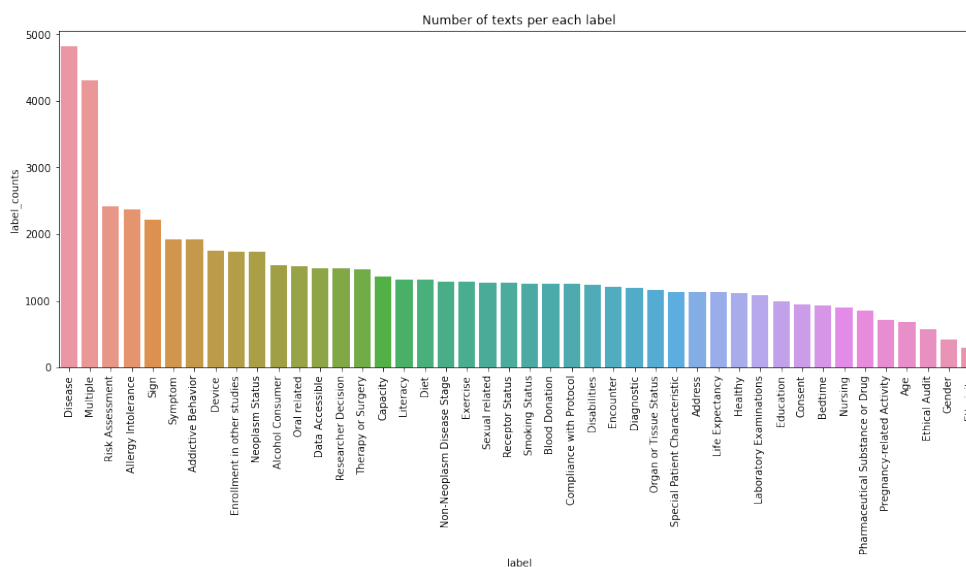


图 2: 数据增强后 label 的数量分布

### 2.3.3 Adapt Language Models to Domains and Tasks [2]

现在的语言模型 (Language Model, LM) 大多是在大量且广泛的文本数据上训练而成的，表现优异。该论文作者思考还有没有必要将模型调整迁移到特定目标任务的领域上 (domain of a target task)，即将已经在

量且广泛文本上预训练过的模型进行第二阶段的预训练，实验表明，不要停止预训练，对于特定的任务，完全可以用任务相关的数据再对语言模型做二次预训练，能大大提高模型性能。

于是，我借鉴该论文的思想，将 CTC 任务的所有数据集一起用于 bert-wwm-ext 的 WWM 预训练任务，再根据二次预训练后的模型在文本分类任务上做 fine-tune，以此验证调整预训练模型的目标域这个方法在该任务上是否有效。

训练超参数设置参考文献 [2]，如表 2 所示：

Computing Infrastructure	RTX 3080ti
Model implementations	<a href="https://github.com/gray311/Bert4textclassification4pl">https://github.com/gray311/Bert4textclassification4pl</a>

Hyperparameter	Assignment
number of steps	100 epochs (TAPT)
batch size	256(through gradient accumulation)
maximum learning rate	1e-4
learning rate optimizer	AdamW
Adam beta weights	(0.9,0.98)
Weight decay	0.01
learning rate decay	Linear

表 2: Hyperparameters for domain- and task- adaptive pretraining

使用 CrossEntropy 作为损失函数调整 WWM 预训练任务，训练过程如图 3 所示。

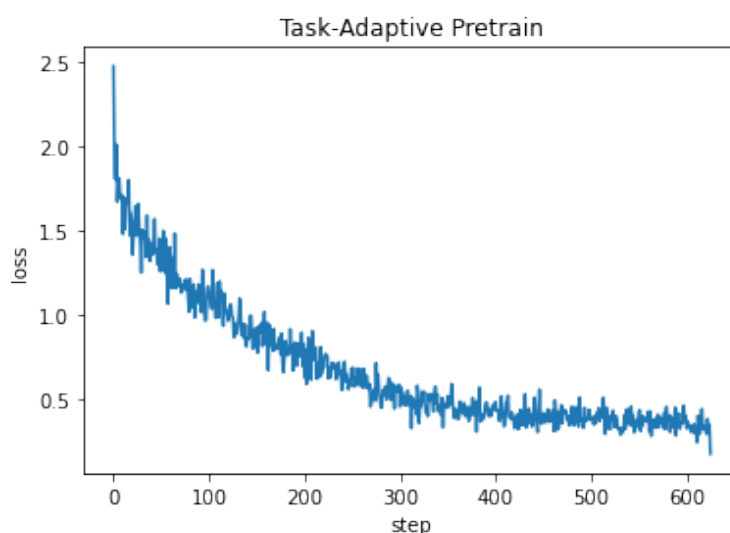


图 3: TAPT 训练过程

## 2.4 Features from different layers [3]

因为 bert 的不同 layer 对于不同下游任务具有不同的适应性 [3]，比如 k 层对 A 任务效果更好，L 层对 B 任务效果更好，所以做法就是尝试用不同层的输出直接完成下游任务，不一定使用 last\_hidden\_layer 作为编码向量。因此，我也尝试了使用 bert-wwm-ext 模型中的不同层输出的特征 token 作为下游任务输出。由于训练

时耗的限制，加上论文中模型错误率表现是随着 layer 层数加深而递减的，因此我只选择了 6-layer,10-layer(0-11layers) 作为特征输出。

### 3 结果

#### 3.1 Hyperparameters for fine-tune

Hyperparameter	Assignment
number of steps	5 epochs
batch size	16
maximum learning rate	2.5e-5
learning rate optimizer	AdamW
Adam beta weights	(0.9,0.999)
Weight decay	None
learning rate decay	None

表 3: Hyperparameters for fine-tune

#### 3.2 benchmark of tricks

本次实验中，每次做 fine-tune 都随机初始 seed，重复 5 次，求得各个指标的平均值，结果如表 3 所示。

Fine-Tune	Dev.			Test		
	Precision	Recall	Macro-F1	Precision	Recall	Macro-F1
baseline	0.7829	0.8543	0.8060	72.385	65.026	66.661
baseline+process	0.7718	0.8712	0.8092	72.066	66.358	67.241
baseline+process+aug	0.7892	0.8523	0.8136	72.339	65.873	67.411
baseline+process+TAPT	<b>0.8135</b>	0.8525	0.8274	72.156	67.412	68.297
baseline+process+aug+TAPT(6layer)	0.7866	0.8773	0.8230	<b>73.615</b>	64.584	67.285
baseline+process+aug+TAPT(10layer)	0.7918	0.8714	0.8254	70.903	<b>68.100</b>	68.136
baseline+process+aug+TAPT(11layer)	0.8006	<b>0.8864</b>	<b>0.8283</b>	73.250	66.389	<b>68.425</b>

表 4: performance with different tricks on CHIP-CTC

将结果上传至天池 CBLUE2.0 榜单进行评测，结果截图如图 4 所示。

1	2022-05-29 21:55:23	baseline+process+TAPT	已完成	4.878	0.0	68.297	72.156	67.412	编辑
2	2022-05-29 21:51:00	baseline+process+aug	已完成	4.815	0.0	67.411	72.339	65.873	编辑
3	2022-05-29 21:47:24	baseline+process	已完成	4.803	0.0	67.241	72.066	66.358	编辑
4	2022-05-29 21:45:02	baseline	已完成	4.761	0.0	66.661	72.385	65.026	编辑

2022-05-29 22:15:54	baseline+process+aug+TAPT(11layer)	已完成	4.887	68.425	73.250	66.389
2022-05-29 22:14:29	baseline+process+aug+TAPT(6layer)	已完成	4.806	67.285	73.615	64.584
2022-05-29 22:05:16	baseline+process+aug+TAPT(10layer)	已完成	4.867	68.136	70.903	68.100

图 4: 天池 CBLUE2.0 榜单评测结果

对比测试结果可以发现, 各项 trick 对模型 Macro-F1 指标都有一定提升效果。在本任务中, bert 模型中最后一层提取的特征效果最好。同时, 考虑到 Recall 指标和 Precision 指标往往不能同时达到最优, 因此可以考虑一些知识蒸馏, 或者集成模型的方法对模型效果进一步提升。另外, 仅靠数据增强生成的 samples 有很大相似性, 所以不能足够有效的解决 label 不平衡的问题, 因此可以考虑用一些 few-sample BERT fine-tuning [4] 方法来提高模型精度。

## 参考文献

- [1] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. Pre-training with whole word masking for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3504–3514, 2021.
- [2] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online, July 2020. Association for Computational Linguistics.
- [3] Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. How to fine-tune BERT for text classification? *CoRR*, abs/1905.05583, 2019.
- [4] Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q. Weinberger, and Yoav Artzi. Revisiting few-sample BERT fine-tuning. *CoRR*, abs/2006.05987, 2020.

## 备注

本文排版基于 L<sup>A</sup>T<sub>E</sub>X [【Github】](#) [【OverLeaf】](#) [【SJTU-ShareLaTeX】](#) 模板, 采用 [【CC BY-NC-SA 4.0】](#) 进行许可。

## A Bert4WWM-Pretrain

```
#!/usr/bin/env python
2 # coding: utf-8

4 # In[1]:

6 import pandas as pd
import torch
8 from datasets import Dataset
import datasets
10 import os
import random
12 import numpy as np

14
def seedeverything(seed):
16     random.seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
18     np.random.seed(seed)
    torch.manual_seed(seed)
20     torch.cuda.manual_seed(seed)
    torch.backends.cudnn.deterministic = True ##
22     torch.backends.cudnn.benchmark = True

24
seedeverything(seed=233)

26
traindata = pd.read_json('./CHIP-CTC/CHIP-CTC_train.json')
28 valdata = pd.read_json('./CHIP-CTC/CHIP-CTC_dev.json')
testdata = pd.read_json('./CHIP-CTC/CHIP-CTC_test.json')
30
examplepreddata = pd.read_excel('./CHIP-CTC/category.xlsx')
32
examplepreddata['label2idx'] = range(examplepreddata.shape[0])
34
label2idx = dict(
36     zip(examplepreddata['Label Name'], examplepreddata['label2idx']))
idx2label = dict(
38     zip(examplepreddata['label2idx'], examplepreddata['Label Name']))

40 print(idx2label)

42 traindata['labels'] = [label2idx[item] for item in traindata['label']]
valdata['labels'] = [label2idx[item] for item in valdata['label']]
44
46 print(len(traindata))
print(len(valdata))
print(len(testdata))
48
traindataset = Dataset.from_pandas(traindata)
50 valdataset = Dataset.from_pandas(valdata)
testdataset = Dataset.from_pandas(testdata)
52
```

```

dataset = datasets.DatasetDict({
54     'train': traindataset,
    'validation': valdataset,
56     'test': testdataset
})

58
print(dataset)

60
train_dataset = dataset['train']
62 val_dataset = dataset['validation']
test_dataset = dataset['test']

64
print(train_dataset.features)

66
print(train_dataset[0])

68
# In[2]:

70
import argparse
72 import json
from typing import List

74
from ltp import LTP
76 from transformers.models.bert.tokenization_bert import BertTokenizer

78
def _is_chinese_char(cp):
80     """Checks whether CP is the codepoint of a CJK character."""
    # This defines a "chinese character" as anything in the CJK Unicode block:
82     # https://en.wikipedia.org/wiki/CJK_Unified_Ideographs_(Unicode_block)
    #
84     # Note that the CJK Unicode block is NOT all Japanese and Korean characters,
    # despite its name. The modern Korean Hangul alphabet is a different block,
86     # as is Japanese Hiragana and Katakana. Those alphabets are used to write
    # space-separated words, so they are not treated specially and handled
88     # like the all of the other languages.
    if ((cp >= 0x4E00 and cp <= 0x9FFF) or (cp >= 0x3400 and cp <= 0x4DBF) #
90         or (cp >= 0x20000 and cp <= 0x2A6DF) #
        or (cp >= 0x2A700 and cp <= 0x2B73F) #
92         or (cp >= 0x2B740 and cp <= 0x2B81F) #
        or (cp >= 0x2B820 and cp <= 0x2CEAF) #
94         or (cp >= 0xF900 and cp <= 0xFAFF)
        or (cp >= 0x2F800 and cp <= 0x2FA1F) #
96     ): #
        return True

98
    return False

100

102 def is_chinese(word: str):
    # word like '180' or '身高' or '神'
104     for char in word:
        char = ord(char)
106         if not _is_chinese_char(char):
            return 0
108     return 1

110

112 def get_chinese_word(tokens: List[str]):
    word_set = set()

114     for token in tokens:

```



```

    chinese_word = len(token) > 1 and is_chinese(token)
116     if chinese_word:
        word_set.add(token)
118 word_list = list(word_set)
    return word_list
120

122 def add_sub_symbol(bert_tokens: List[str], chinese_word_set: set()):
    if not chinese_word_set:
124         return bert_tokens
    max_word_len = max([len(w) for w in chinese_word_set])

126
    bert_word = bert_tokens
    start, end = 0, len(bert_word)
    while start < end:
130         single_word = True
        if is_chinese(bert_word[start]):
132             l = min(end - start, max_word_len)
            for i in range(l, 1, -1):
134                 whole_word = "".join(bert_word[start:start + i])
                if whole_word in chinese_word_set:
136                     for j in range(start + 1, start + i):
                        bert_word[j] = "##" + bert_word[j]
138                     start = start + i
                        single_word = False
140                     break
                if single_word:
142                     start += 1
    return bert_word
144

146 def prepare_ref(lines: List[str], ltp_tokenizer: LTP,
    bert_tokenizer: BertTokenizer):
148     ltp_res = []

150     for i in range(0, len(lines), 100):
        res = ltp_tokenizer.seg(lines[i:i + 100])[0]
152         res = [get_chinese_word(r) for r in res]
        ltp_res.extend(res)
154     assert len(ltp_res) == len(lines)

156     bert_res = []
    for i in range(0, len(lines), 100):
158         res = bert_tokenizer(lines[i:i + 100],
            add_special_tokens=True,
160            truncation=True,
            max_length=512)
        bert_res.extend(res["input_ids"])
162     assert len(bert_res) == len(lines)

164
    ref_ids = []
166     for input_ids, chinese_word in zip(bert_res, ltp_res):

168
        input_tokens = []
        for id in input_ids:
170             token = bert_tokenizer._convert_id_to_token(id)
            input_tokens.append(token)
172         input_tokens = add_sub_symbol(input_tokens, chinese_word)
        ref_id = []
174         # We only save pos of chinese subwords start with ##, which mean is part of a whole word.
        for i, token in enumerate(input_tokens):
176             if token[:2] == "##":

```

```

        clean_token = token[2:]
178         # save chinese tokens' pos
        if len(clean_token) == 1 and _is_chinese_char(
180             ord(clean_token)):
            ref_id.append(i)
182         ref_ids.append(ref_id)

184     assert len(ref_ids) == len(bert_res)

186     return ref_ids

188
# In[3]:
190
from transformers.data.data_collator import DataCollatorForLanguageModeling, DataCollatorForWholeWordMask
192 from transformers import BertForMaskedLM

194 path = "hflchinese-bert-wwm-ext"
tokenizer = BertTokenizer.from_pretrained(path)
196 net = BertForMaskedLM.from_pretrained(path)

198 from ltp import LTP

200 ltp = LTP()

202 sent = [item['text'] for item in train_dataset] + [
        item['text'] for item in val_dataset
204 ] + [item['text'] for item in test_dataset]

206 ref = prepare_ref(sent, ltp, tokenizer)
print(len(ref))
208
# In[4]:
210
from transformers import BertTokenizer, BertModel, AutoModelForSequenceClassification
212 from datasets import Dataset

214 tokenizer = BertTokenizer.from_pretrained("hflchinese-bert-wwm-ext")

216
def tokenize_function(sample):
218     return tokenizer(sample['text'], truncation=True)

220
tokenized_datasets = dataset.map(tokenize_function, batched=True)
222
tokenized_datasets['train'] = tokenized_datasets['train'].remove_columns(
224     ['id', 'text', 'label'])
tokenized_datasets['validation'] = tokenized_datasets[
226     'validation'].remove_columns(['id', 'text', 'label'])
tokenized_datasets['test'] = tokenized_datasets['test'].remove_columns(
228     ['id', 'text'])

230 # In[5]:

232 encoder_dict = tokenized_datasets['train']['input_ids'] + tokenized_datasets[
        'validation']['input_ids'] + tokenized_datasets['test']['input_ids']
234
# 加上子字信息, 而且传入的是List, 不是tensor。
236 train_mlm_dataset = [{
        'input_ids': encoder_dict[i],
238     'chinese_ref': ref[i]

```

```

    } for i in range(len(ref))]

240
datacollector = DataCollatorForWholeWordMask(tokenizer)
242
# In[6]:
244
from torch.utils.data import DataLoader
246 from transformers import DataCollatorWithPadding #实现按batch自动padding

248 train_mlm_dataloader = DataLoader(train_mlm_dataset,
                                   shuffle=True,
250                                   batch_size=8,
                                   collate_fn=datacollector)
252 for batch in train_mlm_dataloader:
    print({k: v.shape for k, v in batch.items()})
254     break

256 # In[7]:

258 for batch in train_mlm_dataloader:
    outputs = net(**batch)
260     print(outputs)
    break
262
# In[8]:
264
import pytorch_lightning as pl
266 from pytorch_lightning import Trainer
from pytorch_lightning.callbacks import EarlyStopping
268 from pytorch_lightning.callbacks import ModelCheckpoint
from tensorboardX import SummaryWriter
270 from transformers import AdamW, get_scheduler

272 from datasets import load_metric
from statistics import mean
274 from sklearn import metrics
from torch import nn
276 import json
import warnings

278 warnings.filterwarnings("ignore")

280
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
282 num_epochs = 30
lr = 2e-5
284
num_training_steps = num_epochs * len(
286     train_mlm_dataloader) # num of batches * num of epochs
print(num_training_steps)
288

290 class Bert4wwmtask_lightningsystem(pl.LightningModule):

292     def __init__(self, net, lr, epoch, len):
        super(Bert4wwmtask_lightningsystem, self).__init__()
294         self.net = net.to(device)
        self.lr = lr
296         self.epoch = epoch
        self.num_training_steps = len
298         self.writer = SummaryWriter('./log')
        self.iteration = 0
300         self.num = 0

```

```

302     #self.metric = load_metric("glue", "mrpc", mirror="tuna")

303
304     def configure_optimizers(self):
305
306         self.optimizer = AdamW(self.net.parameters(),
307                                lr=self.lr,
308                                betas=(0.9, 0.98),
309                                weight_decay=0.01)
310         self.lr_scheduler = get_scheduler(
311             'linear',
312             optimizer=self.optimizer,
313             num_warmup_steps=0,
314             num_training_steps=self.num_training_steps)
315         optim_dict = {
316             'optimizer': self.optimizer,
317             'lr_scheduler': self.lr_scheduler
318         }
319         return optim_dict
320
321     def training_step(self, batch, batch_idx):
322         batch = {k: v.to(device) for k, v in batch.items()}
323         loss = self.net(**batch).loss
324         lr_ = self.lr * (1.0 - self.iteration / self.num_training_steps)**0.9
325         for param_group in self.optimizer.param_groups:
326             param_group['lr'] = lr_
327         self.iteration += 1
328         self.writer.add_scalar(
329             'info/lr',
330             self.optimizer.state_dict()['param_groups'][0]['lr'],
331             self.iteration)
332         self.writer.add_scalar('info/loss', loss, self.iteration)
333         return loss
334
335     def validation_step(self, batch, batch_idx):
336         pass
337
338     def test_step(self, batch, batch_idx):
339         pass
340
341     def training_epoch_end(self, outputs):
342         self.num += 1
343         if self.num % 10 == 0:
344             pt_save_directory = "./pt_save_pretrained" + str(self.num)
345             self.net.save_pretrained(pt_save_directory)
346
347     def validation_epoch_end(self, outputs):
348         pass
349
350 # In[9]:
351
352 model = Bert4wwmtask_lightningsystem(net=net,
353                                     lr=lr,
354                                     epoch=num_epochs,
355                                     len=num_training_steps)
356
357 trainer = Trainer(
358     logger=False,
359     max_epochs=num_epochs,
360     gpus=1,
361     reload_data loaders_every_n_epochs=False,
362     num_sanity_val_steps=0, # Skip Sanity Check

```

```

    #callbacks=[checkpoint_callback],
364     #limit_train_batches=0.05
    precision=16,
366     accumulate_grad_batches=32,
    #gradient_clip_val=0.5,
368 )

370 trainer.fit(model, train_mlm_dataloader)

```

src/bert4wwm.py

## B Bert4CTC-Fine-tune

```

#!/usr/bin/env python
2 # coding: utf-8

4 # In[5]:

6 import pandas as pd
import torch
8 from datasets import Dataset
import datasets
10 import os
import random
12 import numpy as np
import re
14 from copy import deepcopy
import json
16 import pytorch_lightning as pl
from pytorch_lightning import Trainer
18 from pytorch_lightning.callbacks import EarlyStopping
from pytorch_lightning.callbacks import ModelCheckpoint
20 from nlpcda import Similarword

22
def seedeverything(seed):
24     random.seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
26     np.random.seed(seed)
    torch.manual_seed(seed)
28     torch.cuda.manual_seed(seed)
    torch.backends.cudnn.deterministic = True ##
30     torch.backends.cudnn.benchmark = True

32
seedeverything(seed=233)

34
traindata = pd.read_json('./CHIP-CTC/CHIP-CTC_train_aug.json')
36 valdata = pd.read_json('./CHIP-CTC/CHIP-CTC_dev.json')
testdata = pd.read_json('./CHIP-CTC/CHIP-CTC_test.json')
38 testdata_temp = deepcopy(testdata)

40
def textclean(x):
42     a = re.findall('[\u4e00-\u9fa5A-Za-z][\S\s]+', x, re.S)
    a = "".join(a)
44     return a

46

```

```

traindata['text'] = traindata['text'].apply(lambda x: textclean(x))
48 valdata['text'] = valdata['text'].apply(lambda x: textclean(x))
testdata['text'] = testdata['text'].apply(lambda x: textclean(x))
50
examplepreddata = pd.read_excel('./CHIP-CTC/category.xlsx')
52
examplepreddata['label2idx'] = range(examplepreddata.shape[0])
54
label2idx = dict(
56     zip(examplepreddata['Label Name'], examplepreddata['label2idx']))
idx2label = dict(
58     zip(examplepreddata['label2idx'], examplepreddata['Label Name']))

60 with open("idx2label.json", "w", encoding="utf-8") as fp:
    json.dump(idx2label, fp, ensure_ascii=False, indent=4)
62
print(idx2label)
64
traindata['labels'] = [label2idx[item] for item in traindata['label']]
66 valdata['labels'] = [label2idx[item] for item in valdata['label']]

68 print(len(traindata))
print(len(valdata))
70 print(len(testdata))

72 traindataset = Dataset.from_pandas(traindata)
valdataset = Dataset.from_pandas(valdata)
74 testdataset = Dataset.from_pandas(testdata)

76 dataset = datasets.DatasetDict({
    'train': traindataset,
78     'validation': valdataset,
    'test': testdataset
80 })

82 print(dataset)

84 train_dataset = dataset['train']
print(train_dataset.features)
86
print(train_dataset[0])
88
# In[4]:
90
from transformers import BertTokenizer, BertModel, AutoModelForSequenceClassification
92
checkpoint = "hflichinese-bert-wwm-withpretrain-ext"
94 tokenizer = BertTokenizer.from_pretrained(checkpoint)

96
def tokenize_function(sample):
98     return tokenizer(sample['text'], truncation=True)

100
tokenized_datasets = dataset.map(tokenize_function, batched=True)
102
tokenized_datasets['train'] = tokenized_datasets['train'].remove_columns(
104     ['id', 'text', 'label'])
tokenized_datasets['validation'] = tokenized_datasets[
106     'validation'].remove_columns(['id', 'text', 'label'])
tokenized_datasets['test'] = tokenized_datasets['test'].remove_columns(
108     ['id', 'text'])

```

```
110 from transformers import DataCollatorWithPadding #实现按batch自动padding
112 data_collator = DataCollatorWithPadding(tokenizer=tokenizer)
114 print(tokenized_datasets)
116 # In[6]:
118 from torch.utils.data import DataLoader, Dataset
120 train_dataloader = DataLoader(tokenized_datasets['train'],
                                shuffle=True,
122                                batch_size=8,
                                collate_fn=data_collator)
124 val_dataloader = DataLoader(tokenized_datasets['validation'],
                              batch_size=8,
126                              collate_fn=data_collator)
test_dataloader = DataLoader(tokenized_datasets['test'],
                              batch_size=8,
128                              collate_fn=data_collator)
130 for batch in test_dataloader:
    print({k: v.shape for k, v in batch.items()})
132 break
134 # In[7]:
136 from transformers import BertTokenizer, BertModel, AutoModelForSequenceClassification, AutoModel
from transformers import BertForSequenceClassification, BertForMaskedLM
138
net = AutoModelForSequenceClassification.from_pretrained(
140     checkpoint, num_labels=examplepreddata.shape[0])
142 for batch in train_dataloader:
    outputs = net(**batch)
144     print(outputs)
    break
146
148 # In[8]:
149 '''
150 3.8546
151 -1.4115e-01
152 4.1694
153 '''
154 from transformers import AdamW, get_scheduler
156 from datasets import load_metric
from statistics import mean
158 from sklearn import metrics
from torch import nn
160 import json
import warnings
162 from tensorboardX import SummaryWriter
164 warnings.filterwarnings("ignore")
166 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
num_epochs = 5
168 lr = 2.5e-5
num_labels = examplepreddata.shape[0]
170 num_training_steps = num_epochs * len(
```

```

    train_dataloader) # num of batches * num of epochs
172 print(num_training_steps)

174
class Mlp(nn.Module):
176
    def __init__(self,
178         in_features,
            hidden_features=1000,
180         out_features=None,
            act_layer=nn.GELU,
182         drop=0.):
        super().__init__()
184         self.fc1 = nn.Linear(in_features, hidden_features)
            self.act = act_layer()
186         self.fc2 = nn.Linear(hidden_features, out_features)
            self.softmax = nn.Softmax(dim=-1)
188         self.drop = nn.Dropout(drop)

190 def forward(self, x):
    x = self.fc1(x)
192     x = self.act(x)
    x = self.drop(x)
194     x = self.fc2(x)
    x = self.drop(x)
196     return self.softmax(x)

198
class Bert4textclassification_lightningsystem(pl.LightningModule):
200
    def __init__(self, net, lr, epoch, len):
202         super(Bert4textclassification_lightningsystem, self).__init__()
            self.net = net.to(device)
204             self.lr = lr
            self.epoch = epoch
206             self.num_training_steps = len
            self.writer = SummaryWriter('./log-' + checkpoint)
208             self.iteration = 0
            #self.metric = load_metric("glue", "mrpc", mirror="tuna")

210
    def configure_optimizers(self):
212
        self.optimizer = AdamW(self.net.parameters(), lr=self.lr)
214         lr_scheduler = get_scheduler(
            'linear',
216             optimizer=self.optimizer,
            num_warmup_steps=0,
218             num_training_steps=self.num_training_steps)
        optim_dict = {
220             'optimizer': self.optimizer,
            'lr_scheduler': lr_scheduler
222         }
        return optim_dict

224
    def metrics_compute(self, mode, outputs):
226         loss = []
            loss.append(outputs[0][mode + '_loss'])
228             predictions = outputs[0]['predictions']
            labels = outputs[0]['labels']
230             for i in range(1, len(outputs)):
                loss.append(outputs[i][mode + '_loss'])
232             predictions = torch.concat(

```



```

        [predictions, outputs[i]['predictions']], dim=0)
234     labels = torch.concat([labels, outputs[i]['labels']], dim=0)
    loss = torch.tensor(loss)
236     predictions = predictions.cpu().detach().numpy()
    labels = labels.cpu().detach().numpy()
238     return loss, predictions, labels

240 def training_step(self, batch, batch_idx):
    batch = {k: v.to(device) for k, v in batch.items()}
242     outputs = self.net(**batch)
    loss = outputs.loss
244     lr_ = self.lr * (1.0 - self.iteration / self.num_training_steps) ** 0.9
    for param_group in self.optimizer.param_groups:
246         param_group['lr'] = lr_
    logits = outputs.logits
248     predictions = torch.argmax(logits, dim=-1)
    metrics_dict = metrics.classification_report(
250         predictions.cpu().detach().numpy(),
        batch['labels'].cpu().detach().numpy(),
252         digits=4,
        output_dict=True)
254     self.writer.add_scalar('info/loss', loss, self.iteration)
    self.writer.add_scalar('info/weighted_avg',
256         metrics_dict['weighted avg']['f1-score'],
        self.iteration)
258     self.iteration += 1
    return loss

260 def validation_step(self, batch, batch_idx):
    batch = {k: v.to(device) for k, v in batch.items()}
262     outputs = self.net(**batch)
    logits = outputs.logits
264     predictions = torch.argmax(logits, dim=-1)
    metrics_dict = metrics.classification_report(
266         predictions.cpu().detach().numpy(),
        batch['labels'].cpu().detach().numpy(),
268         digits=4,
        output_dict=True)
270     self.log('macro_avg', metrics_dict['macro avg']['f1-score'])
    #self.metric.add_batch(predictions=predictions, references=batch["labels"])
272     return {
        'val_loss': outputs.loss,
        'predictions': predictions,
274         'labels': batch['labels']
    }

278 def test_step(self, batch, batch_idx):
    batch = {k: v.to(device) for k, v in batch.items()}
280     outputs = self.net(**batch)
    logits = outputs.logits
282     predictions = torch.argmax(logits, dim=-1)
    return {'test_loss': outputs.loss, 'predictions': predictions}

284 def training_epoch_end(self, outputs):
    pass

288 def validation_epoch_end(self, outputs):
    print(outputs[0]['predictions'].shape)
    print(len(outputs))
290     val_loss, predictions, labels = self.metrics_compute('val', outputs)
    print(predictions.shape)
292     print('\n', "val_loss: ", val_loss.mean())
294

```

```

print(metrics.classification_report(predictions, labels, digits=4))

296
def test_epoch_end(self, outputs):
298     predictions = outputs[0]['predictions']
    for i in range(1, len(outputs)):
300         predictions = torch.concat(
            [predictions, outputs[i]['predictions']], dim=0)
302     predictions = predictions.cpu().detach().numpy().tolist()
    test_labels = [idx2label[idx] for idx in predictions]
304     testdata_temp['label'] = test_labels
    test_pred_list = []
306     for i in range(testdata_temp.shape[0]):
        temp_dict = {}
308         temp_dict['id'] = testdata_temp.iloc[i, 0]
        temp_dict['label'] = testdata_temp.iloc[i, 2]
310         temp_dict['text'] = testdata_temp.iloc[i, 1]
        test_pred_list.append(temp_dict)
312     print('\n', testdata_temp.head())
    with open("result.json", "w", encoding="utf-8") as fp:
314         json.dump(test_pred_list, fp, ensure_ascii=False, indent=4)

316
# In[ ]:
318
model = Bert4textclassification_lightningsystem(net, lr, num_epochs,
320                                             num_training_steps)
checkpoint_callback = ModelCheckpoint(
322     monitor='macro_avg',
    dirpath='./output_baseline',
324     filename='hflchinese-bert-wwm-ext-CHIP-CTC-{epoch:02d}-{macro_avg:.4f}',
    mode='max',
326     save_top_k=-1,
)
328 trainer = Trainer(
    logger=False,
330     max_epochs=num_epochs,
    gpus=1,
332     reload_dataloaders_every_n_epochs=False,
    num_sanity_val_steps=0, # Skip Sanity Check
334     callbacks=[checkpoint_callback],
    #limit_train_batches=0.05
336     #precision=16,
    accumulate_grad_batches=2,
338     #gradient_clip_val=0.5,
)
340
trainer.fit(model, train_dataloader, val_dataloader)
342
# In[ ]:
344
print(testdata.shape)
346 print(testdata_temp.shape)

348 model = Bert4textclassification_lightningsystem.load_from_checkpoint(
    checkpoint_path=
350     './outputhflchinese-bert-wwm-withpretrain-ext/hflchinese-bert-wwm-ext-CHIP-CTC-epoch=03-weighted_avg=0.8530.ckpt',
    net=net,
352     lr=lr,
    epoch=num_epochs,
354     len=num_training_steps)

356 trainer = Trainer(

```

```

    logger=False,
358   gpus=1,
    #limit_train_batches=0.05
360   #precision=16,
    #accumulate_grad_batches=2,
362   #gradient_clip_val=0.5,
)
364 trainer.test(model=model, dataloaders=test_dataloader)

```

src/bert4CTC.py

## C Chinese Data Augmentation

```

#!/usr/bin/env python
2 # coding: utf-8

4 # In[51]:

6 import pandas as pd
  import torch
8 from datasets import Dataset
  import datasets
10 import os
  import random
12 import numpy as np
  import re
14 from copy import deepcopy
  import json
16 import pytorch_lightning as pl
  from pytorch_lightning import Trainer
18 from pytorch_lightning.callbacks import EarlyStopping
  from pytorch_lightning.callbacks import ModelCheckpoint
20 from nlpcda import Similarword

22
def seedeverything(seed):
24     random.seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
26     np.random.seed(seed)
    torch.manual_seed(seed)
28     torch.cuda.manual_seed(seed)
    torch.backends.cudnn.deterministic = True ##
30     torch.backends.cudnn.benchmark = True

32
seedeverything(seed=233)
34
traindata = pd.read_json('./CHIP-CTC/CHIP-CTC_train.json')
36 valdata = pd.read_json('./CHIP-CTC/CHIP-CTC_dev.json')
  testdata = pd.read_json('./CHIP-CTC/CHIP-CTC_test.json')
38 testdata_temp = deepcopy(testdata)

40 print(traindata.head())

42 # In[52]:

44 import matplotlib.pyplot as plt
  import seaborn as sns
46

```

```

# In[53]:
48 temp = traindata.groupby(["label"])['text'].nunique()
49 df = pd.DataFrame({'label': temp.index, 'label_counts': temp.values})
50 df = df.sort_values(['label_counts'], ascending=False).head(50)
51 plt.figure(figsize=(16, 6))
52 plt.title(f'Number of texts per each label')
53 sns.set_color_codes("pastel")
54 s = sns.barplot(x='label', y="label_counts", data=df)
55 s.set_xticklabels(s.get_xticklabels(), rotation=90)
56 locs, labels = plt.xticks()
57 plt.show()

# In[54]:
60
61 print(df.head())
62 label2num = dict(zip(df['label'], df['label_counts']))
63 print(label2num)
64 print(label2num['Therapy or Surgery'])
65
# In[57]:
66
70 class TextAugmentation:

72     def __init__(self, traindata, valdata, testdata, label2num, padding,
73                  change_rate):
74         super(TextAugmentation, self).__init__()
75         self.traindata = traindata
76         self.valdata = valdata
77         self.testdata = testdata
78         self.label2num = label2num
79         self.padding = padding
80         self.change_rate = change_rate

82     def textclean(self, x):
83         a = re.findall('[\u4e00-\u9fa5A-Za-z][\S\s]+', x, re.S)
84         a = "".join(a)
85         return a

86     def TextProcess(self):
87         self.traindata['text'] = self.traindata['text'].apply(
88             lambda x: self.textclean(x))
89         self.valdata['text'] = self.valdata['text'].apply(
90             lambda x: self.textclean(x))
91         self.testdata['text'] = self.testdata['text'].apply(
92             lambda x: self.textclean(x))

94     def TextAug(self):
95         self.traindata['res'] = ""
96         for i in range(self.traindata.shape[0]):
97             if self.padding > self.label2num[self.traindata.iloc[i, 1]]:
98                 res = max(
99                     math.ceil(
100                         math.sqrt((self.padding -
101                                     self.label2num[self.traindata.iloc[i, 1]]) /
102                                     self.label2num[self.traindata.iloc[i, 1]])),
103                     1)
104             else:
105                 res = 1
106             self.traindata.iloc[i, 3] = res
107
108

```

```

110     for i in tqdm(range(self.traindata.shape[0])):
111         a = self.traindata.iloc[i, 2]
112         res = self.traindata.iloc[i, 3]
113         smw = CharPositionExchange(create_num=res,
114                                   change_rate=self.change_rate,
115                                   char_gram=3)
116         rs1 = smw.replace(a)[1:]
117         if len(rs1) > 0:
118             temp = self.traindata.iloc[i]
119             temp = temp.to_frame()
120             temp = pd.DataFrame(temp.values.T, columns=temp.index)
121             temp = pd.concat([temp] * len(rs1), ignore_index=True)
122             temp['text'] = rs1
123             self.traindata = pd.concat([self.traindata, temp],
124                                       ignore_index=True)
125
126     for i in tqdm(range(self.traindata.shape[0])):
127         a = self.traindata.iloc[i, 2]
128         res = self.traindata.iloc[i, 3]
129         smw = RandomDeleteChar(create_num=res,
130                                change_rate=self.change_rate)
131         rs1 = smw.replace(a)[1:]
132         if len(rs1) > 0:
133             temp = self.traindata.iloc[i]
134             temp = temp.to_frame()
135             temp = pd.DataFrame(temp.values.T, columns=temp.index)
136             temp = pd.concat([temp] * len(rs1), ignore_index=True)
137             temp['text'] = rs1
138             self.traindata = pd.concat([self.traindata, temp],
139                                       ignore_index=True)
140
141     def get_aug_text(self):
142         return self.traindata, self.valdata, self.testdata
143
144 # In[58]:
145
146 temp_train_data = traindata
147 df = temp_train_data
148 temp = df.groupby(["label"])['text'].nunique()
149 df = pd.DataFrame({'label': temp.index, 'label_counts': temp.values})
150 df = df.sort_values(['label_counts'], ascending=False).head(50)
151 plt.figure(figsize=(16, 6))
152 plt.title(f'Number of texts per each label')
153 sns.set_color_codes("pastel")
154 s = sns.barplot(x='label', y="label_counts", data=df)
155 s.set_xticklabels(s.get_xticklabels(), rotation=90)
156 locs, labels = plt.xticks()
157 plt.show()
158
159 # In[60]:
160
161 label2num = dict(zip(df['label'], df['label_counts']))
162 text_aug = TextAugmentation(temp_train_data,
163                             valdata,
164                             testdata,
165                             label2num,
166                             padding=1500,
167                             change_rate=0.3)
168 text_aug.TextProcess()
169 text_aug.TextAug()
170 traindata, valdata, testdata = text_aug.get_aug_text()

```

```
172 print(len(traindata))
    print(len(valdata))
174 print(len(testdata))

176 # In[61]:

178 df = traindata
    print(traindata)
180 temp = df.groupby(["label"])['text'].nunique()
    df = pd.DataFrame({'label': temp.index, 'label_counts': temp.values})
182 df = df.sort_values(['label_counts'], ascending=False).head(50)
    plt.figure(figsize=(16, 6))
184 plt.title(f'Number of texts per each label')
    sns.set_color_codes("pastel")
186 s = sns.barplot(x='label', y="label_counts", data=df)
    s.set_xticklabels(s.get_xticklabels(), rotation=90)
188 locs, labels = plt.xticks()
    plt.show()
```

src/data\_augmentation.py