# Python Challenging项目实验报告

姓名：卢麒萱　学号：2010519

## 项目题目

- 定义：给定一个信息的标题、出处、相关链接以及相关评论，尝试别信息真伪。
- 输入：信息来源、标题、相关超链接、评论
- 输出：真伪标签（0: 消息为真，1: 消息为假）

## 数据分析

1. 数据获取

- https://github.com/yaqingwang/WeFEND-AAAI20
- 随本文件同URL提供
- 只使用有标签数据

2. 数据读取

- 文件格式为csv格式
- 可以使用Python自带的文件读取方式，手动分列
- 可以使用Pandas库进行csv文件读取
- 文件读取代码可以参考上文提及的git仓库中代码

3. 参考读取代码

```
1  with open(filename, 'r', encoding='utf') as f:
2  Import pandas as pd; dataset = pd.read_csv(filename)
```

## 深度学习方法

- PaddlePaddle框架+PaddleNLP

  PaddleNLP和PaddlePaddle框架是什么关系？

Paddle框架是基础底座，提供深度学习任务全流程API。PaddleNLP基于Paddle框架开发，适用于NLP任务

- 使用飞桨完成深度学习任务的通用流程

  - 数据集和数据处理
    paddle.io.Dataset
    paddle.io.DataLoader
    paddlenlp.data
  - 组网和网络配置
    paddle.nn.Embedding
    paddlenlp.seq2vec paddle.nn.Linear
    paddle.tanh paddle.nn.CrossEntropyLoss
    paddle.metric.Accuracy
    paddle.optimizer
    model.prepare
  - 网络训练和评估
    model.fit
    model.evaluate
  - 预测
    model.predict

- 引入相关库

```
1   import os
2   import random
3   import csv
4   import numpy as np
5   from functools import partial
6   import paddle
7   import paddle.nn as nn
8   import paddle.nn.functional as F
9   import paddlenlp as ppnlp
10  from paddlenlp.data import Pad, Stack, Tuple
11  from utils import load_vocab, convert_example
```

- 训练数据、测试数据的文件路径

```
1   original_data = 'challenging/data/train.csv'
2   test_data = 'challenging/data/test.csv'
```

- 解析训练数据

```
1   with open(original_data, "r", encoding='utf-8') as f:
2       text = csv.reader(f)
3       for i in text:
4           if i[5] != "1" and i[5] != "0":
5               continue
6           if i[5] == "1":
7               all_rumor_list.append(i[0] + i[1] + "\t" + i[5] + "\n")
8               rumor_num += 1
9           else:
10              all_non_rumor_list.append(i[0] + i[1] + "\t" + i[5] + "\n")
11              non_rumor_num += 1
12
13  with open(test_data, "r", encoding='utf-8') as f:
14      text = csv.reader(f)
15      for i in text:
16          if i[5] != "1" and i[5] != "0":
17              continue
18          test_list.append(i[0] + i[1] + "\t" + i[5] + "\n")
19          test_num += 1
20
21  print("谣言数据总量为: " + str(rumor_num))
22  print("非谣言数据总量为: " + str(non_rumor_num))
```

- 打乱数据，写入文件

```
1   data_list_path = "data/"
2   all_data_path = data_list_path + "all_data.txt"
3   test_data_path = data_list_path + "test_data.txt"
4   all_data_list = all_rumor_list + all_non_rumor_list
5
6   random.shuffle(all_data_list)
7   random.shuffle(test_list)
8
9   #在生成all_data.txt之前，首先将其清空
10  with open(all_data_path, 'w') as f:
11      f.seek(0)
12      f.truncate()
13
14  with open(test_data_path, 'w') as f:
15      f.seek(0)
16      f.truncate()
17
18  with open(all_data_path, 'a') as f:
19      for data in all_data_list:
20          f.write(data)
21
22  with open(test_data_path, 'a') as f:
23      for data in test_list:
24          f.write(data)
```

- 划分数据集

  创建序列化表示的数据,并按照一定比例划分训练数据与验证数据

```python
with open(os.path.join(data_list_path, 'eval_list.txt'), 'w',
        encoding='utf-8') as f_eval:
    f_eval.seek(0)
    f_eval.truncate()

with open(os.path.join(data_list_path, 'train_list.txt'),
        'w',
        encoding='utf-8') as f_train:
    f_train.seek(0)
    f_train.truncate()

with open(os.path.join(data_list_path, 'all_data.txt'), 'r',
        encoding='utf-8') as f_data:
    lines = f_data.readlines()

with open(os.path.join(data_list_path, 'test_data.txt'), 'r',
        encoding='utf-8') as f_test_init:
    eval_lines = f_test_init.readlines()

i = 0
with open(os.path.join(data_list_path, 'train_list.txt'),
        'a',
        encoding='utf-8') as f_train:
    for line in lines:
        words = line.split('\t')[-1].replace('\n', '')
        label = line.split('\t')[0]
        labs = ""
        labs = label + '\t' + words + '\n'
        f_train.write(labs)
        i += 1

j = 0
with open(os.path.join(data_list_path, 'eval_list.txt'), 'a',
        encoding='utf-8') as f_test:
    for line in eval_lines:
        words = line.split('\t')[-1].replace('\n', '')
        label = line.split('\t')[0]
        labs = ""
        labs = label + '\t' + words + '\n'
        f_test.write(labs)
        j += 1

print("数据列表生成完成! ")
```

- 数据集和数据处理

  自定义数据集

  映射式(map-style)数据集需要继承 `paddle.io.Dataset`

  - `__getitem__`: 根据给定索引获取数据集中指定样本，在 paddle.io.DataLoader 中需要使用此函数通过下标获取样本。
  - `__len__`: 返回数据集样本个数， paddle.io.BatchSampler 中需要样本个数生成下标序列。

```
1  class SelfDefinedDataset(paddle.io.Dataset):
2      def __init__(self, data):
3          super(SelfDefinedDataset, self).__init__()
4          self.data = data
5
6      def __getitem__(self, idx):
7          return self.data[idx]
8
9      def __len__(self):
10          return len(self.data)
11
12      def get_labels(self):
13          return ["0", "1"]
14
15
16  def txt_to_list(file_name):
17      res_list = []
18      for line in open(file_name):
19          res_list.append(line.strip().split('\t'))
20      return res_list
21
22
23  trainlst = txt_to_list('data/train_list.txt')
24  devlst = txt_to_list('data/eval_list.txt')
```

看看数据长什么样

```
1  label_list = train_ds.get_labels()
2  print(label_list)
3
4  for i in range(10):
5      print(train_ds[i])
```

- 数据处理
- 为了将原始数据处理成模型可以读入的格式，本项目将对数据作以下处理：
  - 首先使用jieba切词，之后将jieba切完后的单词映射词表中单词id。



  - 使用 `paddle.io.DataLoader` 接口多线程异步加载数据。

其中用到了PaddleNLP中关于数据处理的API。PaddleNLP提供了许多关于NLP任务中构建有效的数据pipeline的常用API

| API | 简介 |
| --- | --- |
| `paddlenlp.data.Stack` | 堆叠N个具有相同shape的输入数据来构建一个batch，它的输入必须具有相同的shape，输出便是这些输入的堆叠组成的batch数据。 |
| `paddlenlp.data.Pad` | 堆叠N个输入数据来构建一个batch，每个输入数据将会被padding到N个输入数据中最大的长度 |
| `paddlenlp.data.Tuple` | 将多个组batch的函数包装在一起 |

```python
1   import jieba
2
3   dict_path = 'data/dict.txt'
4
5   #创建数据字典，存放位置：dicts.txt。在生成之前先清空dict.txt
6   #在生成all_data.txt之前，首先将其清空
7   with open(dict_path, 'w') as f:
8       f.seek(0)
9       f.truncate()
10
11  dict_set = set()
12  train_data = open('data/train_list.txt')
13  for data in train_data:
14      seg = jieba.lcut(data[:-3])
15      for datas in seg:
16          if not datas is " ":
17              dict_set.add(datas)
18
19  dicts = open(dict_path, 'w')
20  dicts.write('[PAD]\n')
21  dicts.write('[UNK]\n')
22  for data in dict_set:
23      dicts.write(data + '\n')
24  dicts.close()
25  vocab = load_vocab(dict_path)
26
27  for k, v in vocab.items():
28      print(k, v)
29      break
```

- 构造dataloder

  下面的 `create_data_loader` 函数用于创建运行和预测时所需的 `DataLoader` 对象。

  - `paddle.io.DataLoader` 返回一个迭代器，该迭代器根据 `batch_sampler` 指定的顺序迭代返回dataset数据。异步加载数据。
  - `batch_sampler`：DataLoader通过 `batch_sampler` 产生的mini-batch索引列表来 dataset 中索引样本并组成mini-batch
  - `collate_fn`：指定如何将样本列表组合为mini-batch数据。传给它参数需要是一个callable对象，需要实现对组建的batch的处理逻辑，并返回每个batch的数据。在这里传入的是 `prepare_input` 函数，对产生的数据进行pad操作，并返回实际长度等。

```python
1
2   # Reads data and generates mini-batches.
3   def create_dataloader(dataset,
```

```
 4                        trans_function=None,
 5                        mode='train',
 6                        batch_size=1,
 7                        pad_token_id=0,
 8                        batchify_fn=None):
 9      if trans_function:
10          dataset = dataset.apply(trans_function, lazy=True)
11
12      # return_list 数据是否以list形式返回
13      # collate_fn   指定如何将样本列表组合为mini-batch数据。传给它参数需要是一个
   callable对象，需要实现对组建的batch的处理逻辑，并返回每个batch的数据。在这里传入的
   是`prepare_input`函数，对产生的数据进行pad操作，并返回实际长度等。
14      dataloader = paddle.io.DataLoader(dataset,
15                                        return_list=True,
16                                        batch_size=batch_size,
17                                        collate_fn=batchify_fn)
18
19      return dataloader
20
21
22  # python中的偏函数partial，把一个函数的某些参数固定住（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。
23  trans_function = partial(convert_example,
24                           vocab=vocab,
25                           unk_token_id=vocab.get('[UNK]', 1),
26                           is_test=False)
27
28  # 将读入的数据batch化处理，便于模型batch化运算。
29  # batch中的每个句子将会padding到这个batch中的文本最大长度batch_max_seq_len。
30  # 当文本长度大于batch_max_seq时，将会截断到batch_max_seq_len；当文本长度小于
   batch_max_seq时，将会padding补齐到batch_max_seq_len.
31  batchify_fn = lambda samples, fn=Tuple(
32      Pad(axis=0, pad_val=vocab['[PAD]']),  # input_ids
33      Stack(dtype="int64"),  # seq len
34      Stack(dtype="int64")  # label
35  ): [data for data in fn(samples)]
36
37  train_loader = create_dataloader(train_ds,
38                                   trans_function=trans_function,
39                                   batch_size=32,
40                                   mode='train',
41                                   batchify_fn=batchify_fn)
42  dev_loader = create_dataloader(dev_ds,
43                                 trans_function=trans_function,
44                                 batch_size=32,
45                                 mode='validation',
46                                 batchify_fn=batchify_fn)
```
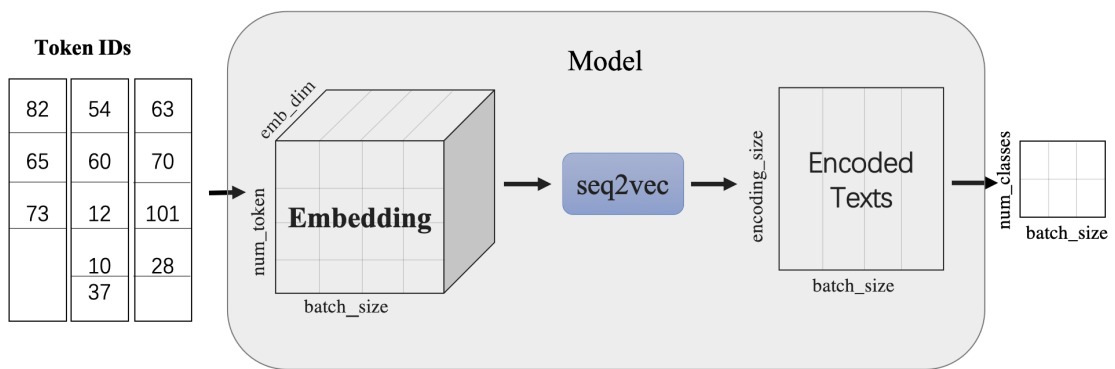
- 模型搭建

  使用 `LSTMencoder` 搭建一个BiLSTM模型用于进行句子建模，得到句子的向量表示。

  然后接一个线性变换层，完成二分类任务。

  - `paddle.nn.Embedding` 组建word-embedding层
  - `ppnlp.seq2vec.LSTMEncoder` 组建句子建模层
  - `paddle.nn.Linear` 构造二分类器

**Token IDs**

| 82 | 54 | 63 |
|----|----|-----|
| 65 | 60 | 70 |
| 73 | 12 | 101 |
|    | 10 | 28 |
|    | 37 |    |

Model

emb_dim · num_token · **Embedding** · batch_size → seq2vec → encoding_size · Encoded Texts · batch_size → num_classes · batch_size

```python
class LSTMModel(nn.Layer):
    def __init__(self,
                 vocab_size,
                 num_classes,
                 emb_dim=128,
                 padding_idx=0,
                 lstm_hidden_size=198,
                 direction='forward',
                 lstm_layers=1,
                 dropout_rate=0,
                 pooling_type=None,
                 fc_hidden_size=96):
        super().__init__()

        # 首先将输入word id 查表后映射成 word embedding
        self.embedder = nn.Embedding(num_embeddings=vocab_size,
                                     embedding_dim=emb_dim,
                                     padding_idx=padding_idx)

        # 将word embedding经过LSTMEncoder变换到文本语义表征空间中
        self.lstm_encoder = ppnlp.seq2vec.LSTMEncoder(
            emb_dim,
            lstm_hidden_size,
            num_layers=lstm_layers,
            direction=direction,
            dropout=dropout_rate,
            pooling_type=pooling_type)

        # LSTMEncoder.get_output_dim()方法可以获取经过encoder之后的文本表示hidden_size
        self.fc = nn.Linear(self.lstm_encoder.get_output_dim(), fc_hidden_size)

        # 最后的分类器
        self.output_layer = nn.Linear(fc_hidden_size, num_classes)

    def forward(self, text, seq_len):
        # text shape: (batch_size, num_tokens)
        # print('input :', text.shape)

        # Shape: (batch_size, num_tokens, embedding_dim)
        embedded_text = self.embedder(text)
        # print('after word-embeding:', embedded_text.shape)
```

```
44          # Shape: (batch_size, num_tokens,
    num_directions*lstm_hidden_size)
45          # num_directions = 2 if direction is 'bidirectional' else 1
46          text_repr = self.lstm_encoder(embedded_text,
    sequence_length=seq_len)
47          # print('after lstm:', text_repr.shape)
48
49          # Shape: (batch_size, fc_hidden_size)
50          fc_out = paddle.tanh(self.fc(text_repr))
51          # print('after Linear classifier:', fc_out.shape)
52
53          # Shape: (batch_size, num_classes)
54          logits = self.output_layer(fc_out)
55          # print('output:', logits.shape)
56
57          # probs 分类概率值
58          probs = F.softmax(logits, axis=-1)
59          # print('output probability:', probs.shape)
60          return probs
61
62
63  model = LSTMModel(len(vocab),
64                    len(label_list),
65                    direction='bidirectional',
66                    padding_idx=vocab['[PAD]'])
67  model = paddle.Model(model)
```

- 模型配置

```
1  optimizer = paddle.optimizer.Adam(parameters=model.parameters(),
2                                    learning_rate=5e-5)
3
4  loss = paddle.nn.CrossEntropyLoss()
5  metric = paddle.metric.Accuracy()
6
7  model.prepare(optimizer, loss, metric)
8  # 设置visualdl路径
9  log_dir = './visualdl'
10 callback = paddle.callbacks.VisualDL(log_dir=log_dir)
```

- 模型训练

  训练过程中会输出loss、acc等信息。这里设置了10个epoch。

```
1  model.fit(train_loader,
2            dev_loader,
3            epochs=10,
4            save_dir='./checkpoints',
5            save_freq=5,
6            callbacks=callback)
```

- 查看训练结果

```
1  results = model.evaluate(dev_loader)
2  print("Finally test acc: %.5f" % results['acc'])
```

- 预测

```
 1  label_map = {1: '谣言', 0: '非谣言'}
 2  results = model.predict(dev_loader, batch_size=128)[0]
 3  predictions = []
 4
 5  for batch_probs in results:
 6      # 映射分类label
 7      idx = np.argmax(batch_probs, axis=-1)
 8      idx = idx.tolist()
 9      labels = [label_map[i] for i in idx]
10      predictions.extend(labels)
11
12  # 看看预测数据前5个样例分类结果
13  for idx, data in enumerate(dev_ds.data[:10]):
14      print('Data: {} \t Label: {}'.format(data[0], predictions[idx]))
```

# 评价指标

- Accuracy



```
Eval samples: 10141
Finally test acc: 0.84311
```



```
Predict begin...
step 317/317 [==============================] - 17ms/step
Predict samples: 10141
Data: Jane心动的信号第六期 | 直男都喜欢什么样的女生? 从刘泽煊回归看各人guan xi    Label: 谣言
Data: 扒爷说【八卦说】某导演找知名网红要L照? 女星离世的玄学论? 模仿前女友上瘾的现女友?         Label: 非谣言
Data: 老唐有态度拿麻将牌"八万"去买车? 某些抖音玩家的低俗令人作呕    Label: 非谣言
Data: 恋爱兮范爷承认了和助理的关系, 李晨: 帽子戴得习惯了, 一波未平一波又起      Label: 谣言
Data: 今日听佛南极冰层现8亿年前女孩?? 进化论已无法解释人类起源! 人从何来? 这里告诉你真相!        Label: 非谣言
Data: 历史奇闻趣事张作霖死前为何阻止一个人进入东北, 历史证明老帅看人非常的准    Label: 非谣言
Data: 我本罪臣Criminal中国环境污染到底有多严重?      Label: 非谣言
Data: 新宁网新宁县人民政府通知: 这些干部已被职务任免!    Label: 非谣言
Data: 小道八卦今天上热搜的某一线女星的老公有不良嗜好!    Label: 非谣言
Data: 江西电影票重磅! 地铁2、3号线最新消息! 昌赣高铁铺轨! 南昌→赣州2小时    Label: 非谣言
```

- Loss



```
step 190/331 - loss: 0.3133 - acc: 0.9838 - 45ms/step
step 200/331 - loss: 0.3446 - acc: 0.9839 - 45ms/step
step 210/331 - loss: 0.3133 - acc: 0.9839 - 45ms/step
step 220/331 - loss: 0.3133 - acc: 0.9844 - 45ms/step
step 230/331 - loss: 0.3133 - acc: 0.9845 - 45ms/step
step 240/331 - loss: 0.3460 - acc: 0.9841 - 45ms/step
step 250/331 - loss: 0.3759 - acc: 0.9840 - 45ms/step
step 260/331 - loss: 0.3445 - acc: 0.9841 - 45ms/step
step 270/331 - loss: 0.3446 - acc: 0.9839 - 45ms/step
step 280/331 - loss: 0.3133 - acc: 0.9837 - 45ms/step
step 290/331 - loss: 0.3758 - acc: 0.9836 - 45ms/step
step 300/331 - loss: 0.3134 - acc: 0.9836 - 45ms/step
step 310/331 - loss: 0.3446 - acc: 0.9831 - 45ms/step
step 320/331 - loss: 0.3133 - acc: 0.9828 - 46ms/step
step 330/331 - loss: 0.3133 - acc: 0.9825 - 46ms/step
step 331/331 - loss: 0.3133 - acc: 0.9825 - 46ms/step
```