homework4

数据集

数据集共包含34282个节点,每个节点为一段文本。除节点属性信息外,还提供了网络的结构信息 (edges.csv)

- train.csv 包含id, name, text, label
- val.csv 包含id, name, text, label
- test.csv 没有label,提交预测文件pred.csv(格式详见submit_example.csv)
- edges.csv 图的结构信息,每行包含一条边连接的两个节点node1,node2
- submit_example.csv 需要提交的预测文件的示例,包含id, label(预测的类别标签)

模型搭建

对于节点数据的处理

利用hugging face 的bert模型进行基于文本的节点分类

在colab上运行

1.挂载数据集

```
from google.colab import drive

drive.mount('/content/drive', force_remount=True)
%cd ..
%cd content/drive/My Drive
```

2.数据读取和查看

```
import pandas as pd

train_data=pd.read_csv("data/train.csv")
val_data=pd.read_csv("data/val.csv")
test_data=pd.read_csv("data/test.csv")
train_data.head()
```

	id	name	text	label
0	21687	"The Perfect Furlough"	1958 films;American films;Comedy films;English	3
1	15781	"Vicky Jenson"	American animators;Living people;Year of birth	0
2	5205	"Billy Burke "	1966 births;American film actors;American tele	0
3	14310	"Giovonnie Samuels"	1985 births;African American actors;American c	0
4	21211	"To Walk with Lions"	1999 films;Canadian films;Kenyan films;English	3

3.应用bert tokenizer并进行fine tuning

bert与原始transformers的主要区别是, BERT没有解码器, 但在基本版本中堆叠了12个编码器。

调用encode_examples函数对训练集以及验证集进行编码并划分batch

tokenizer使用预训练的bert-base-uncased

```
from transformers import BertTokenizer

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
# train dataset
ds_train_encoded =
encode_examples(train_data).shuffle(10000).batch(batch_size)
# val dataset
ds_val_encoded = encode_examples(val_data).batch(batch_size)
```

通过encode_plus将token映射到词嵌入,将文本分词后创建一个包含对应 id, token 类型及是否遮盖的词典。加入特殊token,如[CLS] token将插入序列的开头,[SEP] token位于末尾。

这里设置return attention mask为true,在词典中包含覆盖。

bert 输入格式 [CLS] + tokens + [SEP] + padding

构建tensor flow 数据集

取出每一行的text和label对应的值,利用convert_example_to_feature函数转化为bert的输入 保存通过tokenizer得到的input_ids_list,input_ids是token切片,用于构建将用作模型输入的序列。

attention mask在[0, 1] 中选择的掩码值: 1 表示未屏蔽的标记,0 表示标记的标记。用0作为填充添加的标记。

```
def encode_examples(ds, limit=-1):
    # prepare list, so that we can build up final TensorFlow dataset from
slices.
    input_ids_list = []
    token_type_ids_list = []
    attention_mask_list = []
    label_list = []
    if (limit > 0):
        ds = ds.take(limit)
```

4.模型训练

学习率设为2e-5。

模型初始化,设置分类类别为4。优化器选择adam。

```
model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased',
num_labels=4)
optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate,epsilon=1e-
08, clipnorm=1)
```

利用SparseCategoricalCrossentropy计算损失以及准确率

```
loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy')
model.compile(optimizer=optimizer, loss=loss, metrics=[metric])
```

拟合模型,设置epoch为8

```
bert_history = model.fit(ds_train_encoded, epochs=number_of_epochs,
validation_data=ds_val_encoded)
```

5.训练效果

验证准确率可达98.87%

```
enqing > Desktop > 当代人工智能 > hw4 > assignment > data > 🚦 homework4.ipynb > 🍨 from transformers import TFBertForSequenceClassification
Epoch 2/8
169/169 [==
                         ========] - 259s 2s/step - loss: 0.0586 - accuracy: 0.9847 - val_loss: 0.0499 -
val_accuracy: 0.9858
Epoch 3/8
                      ==============] - 259s 2s/<mark>step</mark> - loss: 0.0443 - accuracy: 0.9873 - val_loss: 0.0520 -
169/169 [=====
val_accuracy: 0.9854
Epoch 4/8
                        val_accuracy: 0.9867
Epoch 5/8
                        -----] - 259s 2s/step - loss: 0.0314 - accuracy: 0.9910 - val_loss: 0.0454 -
169/169 [=====
val_accuracy: 0.9879
Epoch 6/8
169/169 [========
                        ========] - 259s 2s/step - loss: 0.0255 - accuracy: 0.9927 - val_loss: 0.0411 -
val_accuracy: 0.9887
Epoch 7/8
169/169 [===
                         =========] - 259s 2s/step - loss: 0.0205 - accuracy: 0.9943 - val_loss: 0.0491 -
val_accuracy: 0.9883
Epoch 8/8
169/169 [=====
                        =======] - 259s 2s/step - loss: 0.0167 - accuracy: 0.9956 - val_loss: 0.0451 -
val_accuracy: 0.9879
```

6.模型预测

复用之前的模型编码,将函数中的label一栏的传参和存储删掉。

```
def map test example to dict(input ids, attention masks, token type ids):
   return {
      "input ids": input ids,
      "token type ids": token type ids,
      "attention mask": attention masks,
 }
def encode_test_examples(ds, limit=-1):
   input_ids_list = []
   token_type_ids_list = []
   attention mask list = []
   label_list = []
   if (limit > 0):
       ds = ds.take(limit)
   for index, row in ds.iterrows():
       review = row["text"]
          print(review)
          print(label)
       bert input = convert example to feature(review)
        input_ids_list.append(bert_input['input_ids'])
        token type_ids_list.append(bert_input['token_type_ids'])
        attention mask list.append(bert input['attention mask'])
    return tf.data.Dataset.from tensor slices((input ids list,
attention_mask_list, token_type_ids_list)).map(map_test_example_to_dict)
ds_test_encoded = encode_test_examples(test_data).batch(batch_size)
```

```
out=model.predict(ds_test_encoded)
```

out.logits为最后输出的结果矩阵

通过np.argmax取每一行最大元素所在的标签作为这一组的标签值,并存储下来

```
import numpy as np
idx = np.argmax(output, axis=1)
idx.shape
df = pd.DataFrame(idx)
df.to_csv('myfile.csv')
```

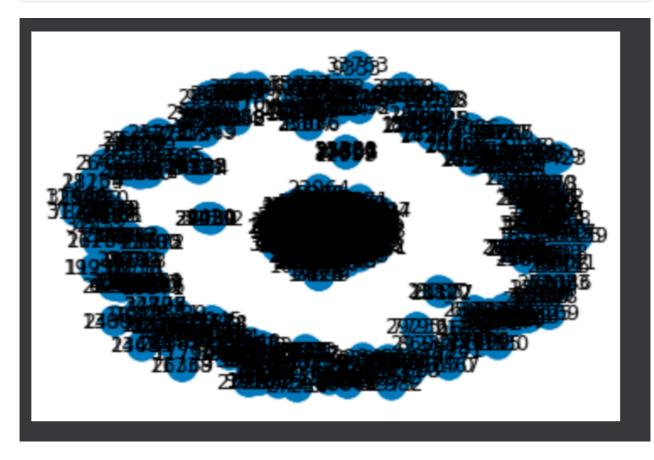
对于图结构信息的研究

尝试画图看一下整体结构,利用networkx库

```
import networkx as nx
import matplotlib.pyplot as plt

for row in edge.itertuples():
    g.add_edge(getattr(row, 'node1'), getattr(row, 'node2'))

fig, ax = plt.subplots()
    nx.draw(g, ax=ax, with_labels=True) # show node label
    plt.show()
```



可以看到网络的连通性较差。

图的基本信息

Node size: 27312 Edge size: 122706

模型效果

单纯使用文本分类效果已经非常好了,由于图的结构非连通,考虑到加入后反而会使得模型训练结果变差,最终使用的模型不包含图结构信息,仅基于文本分类。val_acc约为99%。

错误处理

1.本地tensorflow环境不匹配

```
[Running] python -u "/Users/wangwenqing/Desktop/当代人工智能/hw4/data/homework4.py"

luntimeError: module compiled against API version 0xe but this version of numpy is 0xd

fraceback (most recent call last):

File "/Users/wangwenqing/Desktop/当代人工智能/hw4/data/homework4.py", line 5, in <module>
        import tensorflow as tf

File "/Users/wangwenqing/opt/anaconda3/lib/python3.7/site-packages/tensorflow/__init__.py", line 37, in <module>
        from tensorflow.python.tools import module_util as _module_util

File "/Users/wangwenqing/opt/anaconda3/lib/python3.7/site-packages/tensorflow/python/__init__.py", line 37, in <module>
        from tensorflow.python.eager import context

File "/Users/wangwenqing/opt/anaconda3/lib/python3.7/site-packages/tensorflow/python/eager/context.py", line 35, in <module>
        from tensorflow.python.client import pywrap_tf_session

File "/Users/wangwenqing/opt/anaconda3/lib/python3.7/site-packages/tensorflow/python/client/pywrap_tf_session.py", line 19, in <mod
        from tensorflow.python.client._pywrap_tf_session import *

ImportError: SystemError: <br/>
ImportError: SystemError: <br/>
Sy
```

决定放到colab上训练

2.colab上由于挂载线程连接故障报错

OSError: [Errno 107] Transport endpoint is not connected

从新创建一个副本来再次尝试挂载

报错即可消失

Reference

https://discuss.huggingface.co/t/getting-outputs-of-mode-predict-per-sentence-input/7037 https://huggingface.co/docs/transformers/internal/tokenization_utils

https://swatimeena989.medium.com/bert-text-classification-using-keras-903671e0207d#c685 https://www.jigizhixin.com/articles/2019-05-16-14