



链滴

矩池云 | Tony 老师解读 Kaggle Twitter 情感分析案例

作者: [matpool](#)

原文链接: <https://ld246.com/article/1586485272929>

来源网站: 链滴

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

<p></p>

<p>今天 Tony 老师给大家带来的案例是 Kaggle 上的 Twitter 的情感分析竞赛。在这个案例中，将用预训练的模型 BERT 来完成对整个竞赛的数据分析。</p>

<h2 id="导入需要的库">导入需要的库</h2>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">import numpy as np</span></span><span class="highlight-line"><span class="highlight-cl">import pandas as d</span></span><span class="highlight-line"><span class="highlight-cl">from math import ceil, floor</span></span><span class="highlight-line"><span class="highlight-cl">import tensorflow as tf</span></span><span class="highlight-line"><span class="highlight-cl">import tensorflow.keras.layers as L</span></span><span class="highlight-line"><span class="highlight-cl">from tensorflow.keras.initializers import TruncatedNormal</span></span><span class="highlight-line"><span class="highlight-cl">from sklearn import model_selection</span></span><span class="highlight-line"><span class="highlight-cl">from transformers import BertConfig, TFBertPreTrainedModel, TFBertMainLayer</span></span><span class="highlight-line"><span class="highlight-cl">from tokenizers import BertWordPieceTokenizer</span></span></code></pre>
```

<h2 id="读取并解释数据">读取并解释数据</h2>

<p>在竞赛中，对数据的理解是非常关键的。因此我们首先要做的就是读取数据，然后查看数据的内以及特点。</p>

<p>先用 pandas 来读取 csv 数据，</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">train_df = pd.read_csv('train.csv')</span></span><span class="highlight-line"><span class="highlight-cl">train_df.dropna(inplace=True)</span></span><span class="highlight-line"><span class="highlight-cl">test_df = pd.read_csv('test.csv')</span></span><span class="highlight-line"><span class="highlight-cl">test_df.loc[:, "selected_text"] = test_df.text.values</span></span><span class="highlight-line"><span class="highlight-cl">submission_df = d.read_csv('sample_submission.csv')</span></span></code></pre>
```

<p>再查看下我们的数据的数量，我们一共有 27485 条训练数据，3535 条测试数据，</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">print("train numbers =", train_df.shape)</span></span><span class="highlight-line"><span class="highlight-cl">print("test numbers =", test_df.shape)</span></span></code></pre>
```

<p>紧接着查看训练数据和测试数据前 10 条表单的字段跟数据，表单中包含了一下几个数据字段：</p>

<p>textID: 文本数据记录的唯一 ID; </p>

```

</li>
<li>
<p>text: 原始语句; </p>
</li>
<li>
<p>selected_text: 表示情感的语句; </p>
</li>
<li>
<p>sentiment: 情感类型, neutral 中立, positive 积极, negative 消极; </p>
</li>
</ol>
<p>从数据中我们可以得出, 目标就是根据现有的情感从原本是语句中选出能代表这个情感的语句分。</p>
<pre> <code class="highlight-chroma"> <span class="highlight-line"> <span class="highlight-cl">train_df.head(10)
</span></span> <span class="highlight-line"> <span class="highlight-cl">test_df.head(10)
</span></span></code> </pre>
<h2 id="定义常量">定义常量</h2>
<pre> <code class="highlight-chroma"> <span class="highlight-line"> <span class="highlight-cl"># bert预训练权重跟数据存放的目录
</span></span> <span class="highlight-line"> <span class="highlight-cl">PATH = "./bert-base-uncased/"
</span></span> <span class="highlight-line"> <span class="highlight-cl"># 语句最大长度
</span></span> <span class="highlight-line"> <span class="highlight-cl">MAX_SEQUENCE_LENGTH = 128
</span></span></code> </pre>
<h2 id="载入词向量">载入词向量</h2>
<p>BERT 是依据一个固定的词向量来进行训练的。因此在竞赛中需要先使用 BertWordPieceTokenizer 来加载这些词向量, 其中的 lowercase=True 表示所有的词向量都是小写。设置大小写不敏感可减少模型对资源的占用。</p>
<pre> <code class="highlight-chroma"> <span class="highlight-line"> <span class="highlight-cl">TOKENIZER = BertWordPieceTokenizer(f"{PATH}/vocab.txt", lowercase=True)
</span></span></code> </pre>
<h2 id="定义数据加载器">定义数据加载器</h2>
<p>定义数据预处理函数</p>
<pre> <code class="highlight-chroma"> <span class="highlight-line"> <span class="highlight-cl">def preprocess(tweet, selected_text, sentiment):
</span></span> <span class="highlight-line"> <span class="highlight-cl"># 将被转成byte strng的原始字符串转成utf-8的字符串
</span></span> <span class="highlight-line"> <span class="highlight-cl">    tweet = tweet.
eencode('utf-8')
</span></span> <span class="highlight-line"> <span class="highlight-cl">    selected_text =
elected_text.decode('utf-8')
</span></span> <span class="highlight-line"> <span class="highlight-cl">    sentiment = se
timent.decode('utf-8')
</span></span> <span class="highlight-line"> <span class="highlight-cl">    tweet = " ".join(
tr(tweet).split())
</span></span> <span class="highlight-line"> <span class="highlight-cl">    selected_text =
".join(str(selected_text).split())
</span></span> <span class="highlight-line"> <span class="highlight-cl"># 标记出selected t

```

xt和text共有的单词

```
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_start, idx_end = None, None
</span></span><span class="highlight-line"><span class="highlight-cl"> for index in (i for i, c in enumerate(tweet) if c == selected_text[0]):
</span></span><span class="highlight-line"><span class="highlight-cl"> if tweet[index + len(selected_text)] == selected_text:
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_start = index
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_end = index + len(selected_text)
</span></span><span class="highlight-line"><span class="highlight-cl"> break
</span></span><span class="highlight-line"><span class="highlight-cl"> intersection = [i for i in range(len(tweet)) if i in idx_start:idx_end]
</span></span><span class="highlight-line"><span class="highlight-cl"> if idx_start != None and idx_end != None:
</span></span><span class="highlight-line"><span class="highlight-cl"> for char_idx in range(idx_start, idx_end):
</span></span><span class="highlight-line"><span class="highlight-cl"> intersection[char_idx] = 1
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> # 对原始数据用词量进行编码, 这里会返回原始数据中的词在词向量中的下标
</span></span><span class="highlight-line"><span class="highlight-cl"> # 和原始数据中每个词向量的单词在文中的起始位置跟结束位置
</span></span><span class="highlight-line"><span class="highlight-cl"> enc = TOKENIZER.encode(tweet)
</span></span><span class="highlight-line"><span class="highlight-cl"> input_ids_orig, offsets = enc.ids, enc.offsets
</span></span><span class="highlight-line"><span class="highlight-cl"> target_idx = []
</span></span><span class="highlight-line"><span class="highlight-cl"> for i, (o1, o2) in enumerate(offsets):
</span></span><span class="highlight-line"><span class="highlight-cl"> if sum(intersection[o1:o2]) > 0:
</span></span><span class="highlight-line"><span class="highlight-cl"> target_idx.append(i)
</span></span><span class="highlight-line"><span class="highlight-cl"> target_start = target_idx[0]
</span></span><span class="highlight-line"><span class="highlight-cl"> target_end = target_idx[-1]
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> sentiment_map = {
</span></span><span class="highlight-line"><span class="highlight-cl"> 'positive': 383,
</span></span><span class="highlight-line"><span class="highlight-cl"> 'negative': 497,
</span></span><span class="highlight-line"><span class="highlight-cl"> 'neutral': 869,
</span></span><span class="highlight-line"><span class="highlight-cl"> }
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> # 将情感标签和原句的词向量组合在一起组成我们新的数据
</span></span><span class="highlight-line"><span class="highlight-cl"> input_ids = [101
```

```

+ [sentiment_map[sentiment]] + [102] + input_ids_orig + [102]
</span></span><span class="highlight-line"><span class="highlight-cl"> input_type_ids
[0] * (len(input_ids_orig) + 4)
</span></span><span class="highlight-line"><span class="highlight-cl"> attention_mask
= [1] * (len(input_ids_orig) + 4)
</span></span><span class="highlight-line"><span class="highlight-cl"> offsets = [(0, 0),
0, 0), (0, 0)] + offsets + [(0, 0)]
</span></span><span class="highlight-line"><span class="highlight-cl"> target_start +=
3
</span></span><span class="highlight-line"><span class="highlight-cl"> target_end +=

</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"># 计算需要padding
g的长度, BERT是以固定长度进行输入的, 因此对于不足的我们需要做padding
</span></span><span class="highlight-line"><span class="highlight-cl"> padding_length
= MAX_SEQUENCE_LENGTH - len(input_ids)
</span></span><span class="highlight-line"><span class="highlight-cl"> if padding_leng
h > 0:
</span></span><span class="highlight-line"><span class="highlight-cl"> input_ids = i
put_ids + ([0] * padding_length)
</span></span><span class="highlight-line"><span class="highlight-cl"> attention_ma
k = attention_mask + ([0] * padding_length)
</span></span><span class="highlight-line"><span class="highlight-cl"> input_type_id
= input_type_ids + ([0] * padding_length)
</span></span><span class="highlight-line"><span class="highlight-cl"> offsets = offs
ts + ([0, 0] * padding_length)
</span></span><span class="highlight-line"><span class="highlight-cl"> elif padding_le
ngth < 0:
</span></span><span class="highlight-line"><span class="highlight-cl"> pass
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> return (
</span></span><span class="highlight-line"><span class="highlight-cl"> input_ids, att
ention_mask, input_type_ids, offsets,
</span></span><span class="highlight-line"><span class="highlight-cl"> target_start,
arget_end, tweet, selected_text, sentiment,
</span></span><span class="highlight-line"><span class="highlight-cl"> )
</span></span></code></pre>
<p>定义数据加载器</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">class TweetDataset
(tf.data.Dataset):
</span></span><span class="highlight-line"><span class="highlight-cl"> outputTypes = (
</span></span><span class="highlight-line"><span class="highlight-cl"> tf.dtypes.int3
, tf.dtypes.int32, tf.dtypes.int32,
</span></span><span class="highlight-line"><span class="highlight-cl"> tf.dtypes.int3
, tf.dtypes.float32, tf.dtypes.float32,
</span></span><span class="highlight-line"><span class="highlight-cl"> tf.dtypes.stri
g, tf.dtypes.string, tf.dtypes.string,
</span></span><span class="highlight-line"><span class="highlight-cl"> )
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> outputShapes =
(

```

```

</span></span><span class="highlight-line"><span class="highlight-cl"> (128,), (128,)
(128,),
</span></span><span class="highlight-line"><span class="highlight-cl"> (128, 2), (),
),
</span></span><span class="highlight-line"><span class="highlight-cl"> (), (), (),
</span></span><span class="highlight-line"><span class="highlight-cl"> )
</span></span><span class="highlight-line"><span class="highlight-cl"> def_generator(
</span></span><span class="highlight-line"><span class="highlight-cl"> weet, selected_text, sentiment):
</span></span><span class="highlight-line"><span class="highlight-cl"> for tw, st, se
n zip(tweet, selected_text, sentiment):
</span></span><span class="highlight-line"><span class="highlight-cl"> yield prep
ocess(tw, st, se)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def __new__(cls,
weet, selected_text, sentiment):
</span></span><span class="highlight-line"><span class="highlight-cl"> return tf.data
Dataset.from_generator(
</span></span><span class="highlight-line"><span class="highlight-cl"> cls.genera
or,
</span></span><span class="highlight-line"><span class="highlight-cl"> output_ty
es=cls.outputTypes,
</span></span><span class="highlight-line"><span class="highlight-cl"> output_sh
pes=cls.outputShapes,
</span></span><span class="highlight-line"><span class="highlight-cl"> args=(twe
t, selected_text, sentiment)
</span></span><span class="highlight-line"><span class="highlight-cl"> )
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> @staticmethod
</span></span><span class="highlight-line"><span class="highlight-cl"> def create(dataf
ame, batch_size, shuffle_buffer_size=-1):
</span></span><span class="highlight-line"><span class="highlight-cl"> dataset = Tw
etDataset(
</span></span><span class="highlight-line"><span class="highlight-cl"> dataframe.
ext.values,
</span></span><span class="highlight-line"><span class="highlight-cl"> dataframe.
elected_text.values,
</span></span><span class="highlight-line"><span class="highlight-cl"> dataframe.
entiment.values
</span></span><span class="highlight-line"><span class="highlight-cl"> )
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> dataset = da
aset.cache()
</span></span><span class="highlight-line"><span class="highlight-cl"> if shuffle_buf
er_size != -1:
</span></span><span class="highlight-line"><span class="highlight-cl"> dataset =
ataset.shuffle(shuffle_buffer_size)
</span></span><span class="highlight-line"><span class="highlight-cl"> dataset = da
aset.batch(batch_size)
</span></span><span class="highlight-line"><span class="highlight-cl"> dataset = da
aset.prefetch(tf.data.experimental.AUTOTUNE)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> return dataset

```

```
</span></span></code></pre>
```

<h2 id="定义模型">定义模型</h2>

<p>我们使用 BERT 模型来进行这次竞赛，这里对 BERT 模型做一些简单的介绍。</p>

<p>BERT 的全称是 Bidirectional Encoder Representation from Transformers，即双向 Transformer 的 Encoder，因为 decoder 是不能获要预测的信息的。</p>

<p>模型的主要创新点都在 pre-train 方法上，即用了 Masked LM 和 Next Sentence Prediction 种方法分别捕捉词语和句子级别 representation。</p>

<p>BERT 主要特点如下：</p>

<p>使用了 Transformer 作为算法的主要框架，Trabsformer 能更彻底的捕捉语句中的双向关系；</p>

<p>使用了 Mask Language Model 和 Next Sentence Prediction 的多任务训练目标；</p>

<p>使用更强大的机器训练更大规模的数据，Google 开源了 BERT 模型，我们可以直接使用 BERT 为 Word2Vec 的转换矩阵并高效的将其应用到自己的任务中。</p>

<p>BERT 的本质是在海量的语料基础上，运行自监督学习方法让单词学习得到一个较好的特征表示</p>

<p>在之后特定任务中，可以直接使用 BERT 的特征表示作为该任务的词嵌入特征。所以 BERT 提供是一个供其它任务迁移学习的模型，该模型可以根据任务微调或者固定之后作为特征提取器。</p>

<p>在竞赛中，我们定义了一个 BertModel 类，里面使用 TFBertPreTrainedModel 来进行推理。</p>

<p>BERT 的输出我们保存在 hidden_states 中，然后将这个得到的 hidden_states 结果在加入到 Dense Layer，最后输出我们需要提取的表示情感的文字的起始位置跟结束位置。</p>

<p>这两个位置信息就是我们需要从原文中提取的词向量的位置。</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">
```

```
</span></span><span class="highlight-line"><span class="highlight-cl">class BertModel(T
```

```
</span></span><span class="highlight-line"><span class="highlight-cl">
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"># drop out rate,
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> dr = 0.1
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"># hidden state数量
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> hs = 2
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> def __init__(self,
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> super().__init__
```

```
</span></span><span class="highlight-line"><span class="highlight-cl">
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> self.bert = TF
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> ertMainLayer(config, name="bert")
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> self.concat =
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> .Concatenate()
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> self.dropout
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> L.Dropout(self.dr)
```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> self.qa_outpu
```

```
s = L.Dense(
```

```

</span></span><span class="highlight-line"><span class="highlight-cl">         config.nu
_labels,
</span></span><span class="highlight-line"><span class="highlight-cl">         kernel_initi
lizer=TruncatedNormal(stddev=config.initializer_range),
</span></span><span class="highlight-line"><span class="highlight-cl">         dtype='flo
t32',
</span></span><span class="highlight-line"><span class="highlight-cl">         name="qa
outputs")
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> @tf.function
</span></span><span class="highlight-line"><span class="highlight-cl"> def call(self, inp
ts, **kwargs):
</span></span><span class="highlight-line"><span class="highlight-cl">         _, _, hidden_s
ates = self.bert(inputs, **kwargs)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">         hidden_states
= self.concat([
</span></span><span class="highlight-line"><span class="highlight-cl">         hidden_sta
es[-i] for i in range(1, self.hs+1)
</span></span><span class="highlight-line"><span class="highlight-cl">     ])
</span></span><span class="highlight-line"><span class="highlight-cl">         hidden_states
= self.dropout(hidden_states, training=kwargs.get("training", False))
</span></span><span class="highlight-line"><span class="highlight-cl">         logits = self.
a_outputs(hidden_states)
</span></span><span class="highlight-line"><span class="highlight-cl">         start_logits,
end_logits = tf.split(logits, 2, axis=-1)
</span></span><span class="highlight-line"><span class="highlight-cl">         start_logits =
tf.squeeze(start_logits, axis=-1)
</span></span><span class="highlight-line"><span class="highlight-cl">         end_logits =
tf.squeeze(end_logits, axis=-1)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">         return start_l
ogits, end_logits
</span></span></code></pre>
<p>定义训练函数</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">def train(model, d
taset, loss_fn, optimizer):
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> @tf.function
</span></span><span class="highlight-line"><span class="highlight-cl"> def train_step(
odel, inputs, y_true, loss_fn, optimizer):
</span></span><span class="highlight-line"><span class="highlight-cl">         with tf.Gradi
ntTape() as tape:
</span></span><span class="highlight-line"><span class="highlight-cl">         y_pred =
odel(inputs, training=True)
</span></span><span class="highlight-line"><span class="highlight-cl">         loss = loss
fn(y_true[0], y_pred[0])
</span></span><span class="highlight-line"><span class="highlight-cl">         loss += los
s_fn(y_true[1], y_pred[1])
</span></span><span class="highlight-line"><span class="highlight-cl">         scaled_loss
= optimizer.get_scaled_loss(loss)
</span></span>
</code></pre>

```



```

</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> scaled_gradi
nts = tape.gradient(scaled_loss, model.trainable_variables)
</span></span><span class="highlight-line"><span class="highlight-cl"> gradients =
optimizer.get_unscaled_gradients(scaled_gradients)
</span></span><span class="highlight-line"><span class="highlight-cl"> optimizer.app
y_gradients(zip(gradients, model.trainable_variables))
</span></span><span class="highlight-line"><span class="highlight-cl"> return loss, y
pred
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> epoch_loss = 0.
</span></span><span class="highlight-line"><span class="highlight-cl"> for batch_num,
ample in enumerate(dataset):
</span></span><span class="highlight-line"><span class="highlight-cl"> loss, y_pred =
train_step(model, sample[:3], sample[4:6], loss_fn, optimizer)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> epoch_loss
= loss
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> print(
</span></span><span class="highlight-line"><span class="highlight-cl"> f"training ..
batch {batch_num+1:03d} : "
</span></span><span class="highlight-line"><span class="highlight-cl"> f"train loss
epoch_loss/(batch_num+1):.3f} ",
</span></span><span class="highlight-line"><span class="highlight-cl"> end='\r')
</span></span></code></pre>
<p>定义预制函数</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">def predict(model, dataset, loss_fn, optimizer):
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> @tf.function
</span></span><span class="highlight-line"><span class="highlight-cl"> def predict_step
model, inputs):
</span></span><span class="highlight-line"><span class="highlight-cl"> return model(
nputs)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def to_numpy(*
rgs):
</span></span><span class="highlight-line"><span class="highlight-cl"> out = []
</span></span><span class="highlight-line"><span class="highlight-cl"> for arg in arg
:
</span></span><span class="highlight-line"><span class="highlight-cl"> if arg.dtyp
== tf.string:
</span></span><span class="highlight-line"><span class="highlight-cl"> arg = [s
decode('utf-8') for s in arg.numpy()]
</span></span><span class="highlight-line"><span class="highlight-cl"> out.app
nd(arg)
</span></span><span class="highlight-line"><span class="highlight-cl"> else:
</span></span><span class="highlight-line"><span class="highlight-cl"> arg = ar
.numpy()
</span></span><span class="highlight-line"><span class="highlight-cl"> out.app
nd(arg)
</span></span><span class="highlight-line"><span class="highlight-cl"> return out
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span></code></pre>

```

```

</span></span><span class="highlight-line"><span class="highlight-cl"> offset = tf.zeros
[0, 128, 2], dtype=tf.dtypes.int32)
</span></span><span class="highlight-line"><span class="highlight-cl"> text = tf.zeros([
,], dtype=tf.dtypes.string)
</span></span><span class="highlight-line"><span class="highlight-cl"> selected_text =
f.zeros([0,], dtype=tf.dtypes.string)
</span></span><span class="highlight-line"><span class="highlight-cl"> sentiment = tf.
eros([0,], dtype=tf.dtypes.string)
</span></span><span class="highlight-line"><span class="highlight-cl"> pred_start = tf.
eros([0, 128], dtype=tf.dtypes.float32)
</span></span><span class="highlight-line"><span class="highlight-cl"> pred_end = tf.z
ros([0, 128], dtype=tf.dtypes.float32)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> for batch_num,
ample in enumerate(dataset):
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> print(f"predic
ing ... batch {batch_num+1:03d}"+" "*20, end='\r')
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> y_pred = pre
dict_step(model, sample[:3])
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> # add batch
o accumulators
</span></span><span class="highlight-line"><span class="highlight-cl"> pred_start = t
.concat((pred_start, y_pred[0]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl"> pred_end = tf
.concat((pred_end, y_pred[1]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl"> offset = tf.co
cat((offset, sample[3]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl"> text = tf.conc
t((text, sample[6]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl"> selected_text
= tf.concat((selected_text, sample[7]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl"> sentiment = t
.concat((sentiment, sample[8]), axis=0)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> pred_start, pred
end, text, selected_text, sentiment, offset = \
</span></span><span class="highlight-line"><span class="highlight-cl"> to_numpy(pr
d_start, pred_end, text, selected_text, sentiment, offset)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> return pred_star
, pred_end, text, selected_text, sentiment, offset
</span></span></code></pre>

```

<p>判断函数</p>

<p>这个竞赛采用单词级 Jaccard 系数，计算公式如下</p>

<p></p>

<p>Jaccard 系数计算的是你预测的单词在数据集中的个数，</p>

```

<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">def jaccard(str1, str2):
</span></span><span class="highlight-line"><span class="highlight-cl"> a = set(str1.low

```

```

r().split())
</span></span><span class="highlight-line"><span class="highlight-cl"> b = set(str2.low
r().split())
</span></span><span class="highlight-line"><span class="highlight-cl"> c = a.intersecti
n(b)
</span></span><span class="highlight-line"><span class="highlight-cl"> return float(len(
)) / (len(a) + len(b) - len(c))
</span></span></code></pre>
<p>定义预测结果解码函数</p>
<p>解码函数通过模型预测拿到的 start 和 end 的 index 位置信息，然后和之前拿到的词向量在样
句子中的位置进行比较，将这个区间内的所有的单词都提取出来作为我们的预测结果。</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">def decode_predic
ion(pred_start, pred_end, text, offset, sentiment):
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def decode(pre
_start, pred_end, text, offset):
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded_text
= ""
</span></span><span class="highlight-line"><span class="highlight-cl"> for i in range
pred_start, pred_end+1):
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded_t
xt += text[offset[i][0]:offset[i][1]]
</span></span><span class="highlight-line"><span class="highlight-cl"> if (i+1) &lt;
len(offset) and offset[i][1] &lt; offset[i+1][0]:
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded
_text += " "
</span></span><span class="highlight-line"><span class="highlight-cl"> return decod
d_text
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded_predic
ions = []
</span></span><span class="highlight-line"><span class="highlight-cl"> for i in range(le
n(text)):
</span></span><span class="highlight-line"><span class="highlight-cl"> if sentiment[i]
== "neutral" or len(text[i].split()) &lt; 2:
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded_t
xt = text[i]
</span></span><span class="highlight-line"><span class="highlight-cl"> else:
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_start =
np.argmax(pred_start[i])
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_end =
p.argmax(pred_end[i])
</span></span><span class="highlight-line"><span class="highlight-cl"> if idx_start
&gt; idx_end:
</span></span><span class="highlight-line"><span class="highlight-cl"> idx_end
= idx_start
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded_t
xt = str(decode(idx_start, idx_end, text[i], offset[i]))
</span></span><span class="highlight-line"><span class="highlight-cl"> if len(deco
ded_text) == 0:
</span></span><span class="highlight-line"><span class="highlight-cl"> decoded

```

```

_text = text[i]
</span></span><span class="highlight-line"><span class="highlight-cl">    decoded_pre
ictions.append(decoded_text)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">    return decoded
predictions
</span></span></code></pre>
<h2 id="开始训练">开始训练</h2>
<p>将训练数据分成 5 个 folds, 每个 folds 训练 5 个 epoch, 使用 adam 优化器, learning rate
置成 3e-5, batch size 使用 32。</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">num_folds = 5
</span></span><span class="highlight-line"><span class="highlight-cl">num_epochs = 5
</span></span><span class="highlight-line"><span class="highlight-cl">batch_size = 32
</span></span><span class="highlight-line"><span class="highlight-cl">learning_rate = 3e
5
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">optimizer = tf.ker
s.optimizers.Adam(learning_rate)
</span></span><span class="highlight-line"><span class="highlight-cl">optimizer = tf.ker
s.mixed_precision.experimental.LossScaleOptimizer(
</span></span><span class="highlight-line"><span class="highlight-cl">    optimizer, 'dyn
mic')
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">config = BertConf
g(output_hidden_states=True, num_labels=2)
</span></span><span class="highlight-line"><span class="highlight-cl">model = BertMod
l.from_pretrained(PATH, config=config)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">loss_fn = tf.keras.l
sses.SparseCategoricalCrossentropy(from_logits=True)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">kfold = model_sel
ction.KFold(
</span></span><span class="highlight-line"><span class="highlight-cl">    n_splits=num_f
lds, shuffle=True, random_state=42)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">test_preds_start =
np.zeros((len(test_df), 128), dtype=np.float32)
</span></span><span class="highlight-line"><span class="highlight-cl">test_preds_end =
p.zeros((len(test_df), 128), dtype=np.float32)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">for fold_num, (trai
_idx, valid_idx) in enumerate(kfold.split(train_df.text)):
</span></span><span class="highlight-line"><span class="highlight-cl">    print("\nfold %
2d" % (fold_num+1))
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"># 创建train, valid,
est数据集
</span></span><span class="highlight-line"><span class="highlight-cl">    train_dataset =
weetDataset.create(
</span></span><span class="highlight-line"><span class="highlight-cl">        train_df.iloc[t
ain_idx], batch_size, shuffle_buffer_size=2048)

```

```

</span></span><span class="highlight-line"><span class="highlight-cl"> valid_dataset =
TweetDataset.create(
</span></span><span class="highlight-line"><span class="highlight-cl">     train_df.iloc[
alid_idx], batch_size, shuffle_buffer_size=-1)
</span></span><span class="highlight-line"><span class="highlight-cl"> test_dataset =
weetDataset.create(
</span></span><span class="highlight-line"><span class="highlight-cl">     test_df, batch
size, shuffle_buffer_size=-1)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> best_score = fl
at('-inf')
</span></span><span class="highlight-line"><span class="highlight-cl"> for epoch_num
n range(num_epochs):
</span></span><span class="highlight-line"><span class="highlight-cl">     print("\nepo
h %03d" % (epoch_num+1))
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     train(model, t
ain_dataset, loss_fn, optimizer)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     pred_start, p
red_end, text, selected_text, sentiment, offset = \
</span></span><span class="highlight-line"><span class="highlight-cl">         predict(m
del, valid_dataset, loss_fn, optimizer)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     selected_text
pred = decode_prediction(
</span></span><span class="highlight-line"><span class="highlight-cl">         pred_start,
pred_end, text, offset, sentiment)
</span></span><span class="highlight-line"><span class="highlight-cl">     jaccards = []
</span></span><span class="highlight-line"><span class="highlight-cl">     for i in range(
len(selected_text)):
</span></span><span class="highlight-line"><span class="highlight-cl">         jaccards.a
ppend(
</span></span><span class="highlight-line"><span class="highlight-cl">             jaccard(
elected_text[i], selected_text_pred[i]))
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     score = np.m
ax(jaccards)
</span></span><span class="highlight-line"><span class="highlight-cl">     print(f"valid j
ccard epoch {epoch_num+1:03d}: {score}"+" "*15)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     if score >
est_score:
</span></span><span class="highlight-line"><span class="highlight-cl">         best_score
= score
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> # predict test set
</span></span><span class="highlight-line"><span class="highlight-cl">     test_pred_s
tart, test_pred_end, test_text, _, test_sentiment, test_offset = \
</span></span><span class="highlight-line"><span class="highlight-cl">         predict
model, test_dataset, loss_fn, optimizer)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">     test_preds_start
+= test_pred_start * 0.2

```

```
</span></span><span class="highlight-line"><span class="highlight-cl"> test_preds_end
+= test_pred_end * 0.2
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"># 重置模型, 避免
OM
</span></span><span class="highlight-line"><span class="highlight-cl"> session = tf.co
pat.v1.get_default_session()
</span></span><span class="highlight-line"><span class="highlight-cl"> graph = tf.com
at.v1.get_default_graph()
</span></span><span class="highlight-line"><span class="highlight-cl"> del session, gra
h, model
</span></span><span class="highlight-line"><span class="highlight-cl"> model = BertM
del.from_pretrained(PATH, config=config)
</span></span></code></pre>
<h2 id="预测测试数据-并生成提交文件">预测测试数据, 并生成提交文件</h2>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">selected_text_pred = decode_prediction(
</span></span><span class="highlight-line"><span class="highlight-cl"> test_preds_start,
test_preds_end, test_text, test_offset, test_sentiment)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">def f(selected):
</span></span><span class="highlight-line"><span class="highlight-cl"> return " ".join(s
t(selected.lower().split()))
</span></span><span class="highlight-line"><span class="highlight-cl">submission_df.loc[
:selected_text'] = selected_text_pred
</span></span><span class="highlight-line"><span class="highlight-cl">submission_df['sel
cted_text'] = submission_df['selected_text'].map(f)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">submission_df.to_
sv("submission.csv", index=False)
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span></code></pre>
<p>这个方案在提交的时候在 553 个队伍中排名 153 位, 分数为 0.68。</p>
<p></p>
<p>Twitter 情感分析案例之后会在<a href="https://ld246.com/forward?goto=https%3A%2F%2
www.matpool.com%2F" target="_blank" rel="nofollow ugc">矩池云</a> Demo 镜像中上线,
以直接使用。另矩池云还支持了 Paddle、MindSpore、MegEngine、Jittor 等国产深度学习框架,
免安装直接运行。</p>
```