

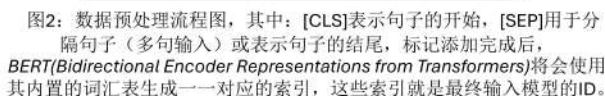
刘沛卓
西北工业大学附属中学
oziul@berkeley.edu 181-8920-2653

本研究使用了Sentiment140数据集，该数据集由Twitter上的用户推文组成，包含了大约160万条数据。

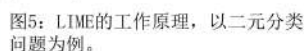
数据包括目标(target)、索引(ids)、日期(date)、标记(flag)、用户(user)、文本(text)共六方面。其中，目标(target)部分表示该推文被标注为积极或消极的情绪，该数据集庞大的数据量和准确的标记使得其成为一个较为理想的训练数据集。



我们将数据集分成了互不重叠的训练集和测试集两部分。此外，我们还对数据进行了预处理，去除了特殊字符，例如标点符号和网址链接，将所有字母均转换为小写，并去除了停用词，以尽量减小以上因素对模型的干扰，并且通过tokenization和添加特殊标记（如[CLS]和[SEP]）将推文文本转换为模型可接受的输入格式。



为了解决深度学习模型的“黑箱”问题[2],本研究引入了LIME(Local Interpretable Model-agnostic Explanations)技术,从而帮助用户理解模型的决策过程。LIME通过提供特征的重要性分析,提升了模型的可解释性,帮助用户进一步相信模型并合理判断是否使用模型的决策。



局部解释模型, 用于解释 f 在 x 附近的行为

$\underset{g \in G}{\operatorname{argmin}}$ 从候选解释模型集合 G 中选择最优模型 g

$L(g, f, \pi_x(z'))$ 损失函数, 衡量简单模型 g 对复杂模型 f 在点 x 邻域内的拟合效果

$\Omega(g)$ 正则化项, 限制模型 g 的复杂性

我们采用最简单的控制台交互界面，用户输入想要分析的语句后，模型会输出情感预测，置信程度和LIME解释：

在LIME解释部分，Feature作为输入特征，代表原始模型做出这个决策的判断依据，而Weight则代表该特征对于决策的贡献程度。

图6: 用户交互界面

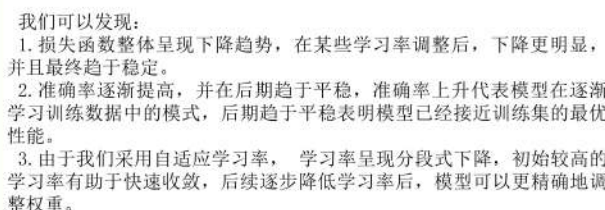
我们调取了Hugging Face的Transformers库，对预训练好的BERT模型进行Fine-tuning，并根据二元分类任务设置num_labels=2，并使用稀疏交叉熵(Sparse Categorical Cross Entropy)损失函数进行模型编译，分类准确率(Precision)作为模型的评估标准(Metrics)。训练过程中，我们使用Adam(Adaptive Moment Estimation)优化器并采用自适应学习率，将批次大小设置为16，并设定4轮训练，在训练过程中通过验证集评估模型表现，防止过拟合。



图4（右）：BERT的实现原理，即Transformers编码器

训练过程中，我们以batch为单位，绘制Loss & Accuracy 曲线，以及学习率的变化情况如下：

Loss: Accuracy: 学习率:



通过K-fold交叉验证的方法，本研究对模型在Sentiment140数据集上的性能进行了评估。具体结果如下：

准确率 (Precision): 对于积极和消极推文的平均准确率分别约为 **96.51%** 和 **95.34%**。

召回率 (Recall): 对于积极和消极推文, 平均召回率分别约为 95.30% 和 96.55%。

F1-score: 对于积极和消极推文的F1-score约为**95.90%**和**95.94%**。

混淆矩阵表示如下:

结果热力图如下:

