

Hugging Face 模型微调训练（基于 BERT 的中文评价情感分析）

目录

- [1. 模型微调的基本概念与流程](#)
- [2. 加载数据集](#)
 - [数据集格式](#)
 - [数据集信息](#)
- [3. 制作 Dataset](#)
 - [数据集字段](#)
 - [数据集信息](#)
- [4. vocab 字典操作](#)
 - [词汇表 \(vocab\)](#)
 - [文本转换](#)
- [5. 下游任务模型设计](#)
 - [模型结构](#)
 - [模型初始化](#)
- [6. 自定义模型训练](#)
 - [数据加载](#)
 - [优化器](#)
 - [训练循环](#)
- [7. 最终效果评估与测试](#)
 - [准确率 \(Accuracy\)](#)
 - [精确率、召回率和 F1 分数](#)
 - [结果分析与模型优化](#)
 - [保存与加载模型](#)
- [8. 课件小结](#)

1. 模型微调的基本概念与流程

微调是指在预训练模型的基础上，通过进一步的训练来适应特定的下游任务。BERT 模型通过预训练来学习语言的通用模式，然后通过微调来适应特定任务，如情感分析、命名实体识别等。微调过程中，通常冻结 BERT 的预训练层，只训练与下游任务相关的层。本课件将介绍如何使用 BERT 模型进行情感分析任务的微调训练。

2. 加载数据集

情感分析任务的数据通常包括文本及其对应的情感标签。使用 Hugging Face 的 datasets 库可以轻松地加载和处理数据集。

```
from datasets import load_dataset
# 加载数据集
dataset = load_dataset('csv', data_files="data/ChnSentiCorp.csv")
# 查看数据集信息
print(dataset)
```

2.1 数据集格式

Hugging Face 的 datasets 库支持多种数据集格式，如 CSV、JSON、TFRecord 等。在本案例中，使用 CSV 格式，CSV 文件应包含两列：一列是文本数据，另一列是情感标签。

2.2 数据集信息

加载数据集后，可以查看数据集的基本信息，如数据集大小、字段名称等。这有助于我们了解数据的分布情况，并在后续步骤中进行适当的处理。

3. 制作 Dataset

加载数据集后，需要对其进行处理以适应模型的输入格式。这包括数据清洗、格式转换等操作。

```
from datasets import Dataset
# 制作 Dataset
dataset = Dataset.from_dict({
    'text': ['位置尚可，但距离海边的位置比预期的要差的多', '5月8日付款成功，当当网显示5月10日发货，可是至今还没看到货物，也没收到任何通知，简不知怎么说好!!!', '整体来说，本书还是不错的。至少在书中描述了许多现实中存在的司法系统方面的问题，这是值得每个法律工作者去思考的。尤其是让那些涉世不深的想加入到律师队伍中的年青人，看到了社会特别是中国司法界真实的一面。缺点是：书中引用了大量的法律条文和司法解释，对于已经是律师或有一定工作经验的法律工作者来说有点多余，而且所占的篇幅不少，有凑字数的嫌疑。整体来说还是不错的。不要对一本书提太高的要求。'],
    'label': [0, 1, 1] # 0 表示负向评价，1 表示正向评价
})
# 查看数据集信息
print(dataset)
```

3.1 数据集字段

在制作 Dataset 时，需定义数据集的字段。在本案例中，定义了两个字段：text（文本）和 label（情感标签）。每个字段都需要与模型的输入和输出匹配。

3.2 数据集信息

制作 Dataset 后，可以通过 dataset.info 等方法查看其大小、字段名称等信息，以确保数据集的正确性和完整性。

4. vocab 字典操作

在微调 BERT 模型之前，需要将模型的词汇表（vocab）与数据集中的文本匹配。这一步骤确保输入的文本能够被正确转换为模型的输入格式。

```
from transformers import BertTokenizer
# 加载 BERT 模型的 vocab 字典
tokenizer = BertTokenizer.from_pretrained('bert-base-chinese')
# 将数据集中的文本转换为 BERT 模型所需的输入格式
dataset = dataset.map(lambda x: tokenizer(x['text'], return_tensors="pt"),
    batched=True)
# 查看数据集信息
print(dataset)
```

4.1 词汇表 (vocab)

BERT 模型使用词汇表 (vocab) 将文本转换为模型可以理解的输入格式。词汇表包含所有模型已知的单词及其对应的索引。确保数据集中的所有文本都能找到对应的词汇索引是至关重要的。

4.2 文本转换

使用 tokenizer 将文本分割成词汇表中的单词，并转换为相应的索引。此步骤需要确保文本长度、特殊字符处理等都与 BERT 模型的预训练设置相一致。

5. 下游任务模型设计

在微调 BERT 模型之前，需要设计一个适应情感分析任务的下游模型结构。通常包括一个或多个全连接层，用于将 BERT 输出的特征向量转换为分类结果。

```
from transformers import BertModel
import torch.nn as nn

class SentimentAnalysisModel(nn.Module):
    def __init__(self):
        super().__init__()
        self.bert = BertModel.from_pretrained('bert-base-chinese')
        self.drop_out = nn.Dropout(0.3)
        self.linear = nn.Linear(768, 2) # 假设情感分类为二分类

    def forward(self, input_ids, attention_mask):
        _, pooled_output = self.bert(
            input_ids=input_ids,
            attention_mask=attention_mask,
            return_dict=False
        )
        output = self.drop_out(pooled_output)
        return self.linear(output)
```

5.1 模型结构

下游任务模型通常包括以下几个部分：

- **BERT 模型**：用于生成文本的上下文特征向量。
- **Dropout 层**：用于防止过拟合，通过随机丢弃一部分神经元来提高模型的泛化能力。
- **全连接层**：用于将 BERT 的输出特征向量映射到具体的分类任务上。

5.2 模型初始化

使用 `BertModel.from_pretrained()` 方法加载预训练的 BERT 模型，同时也可以初始化自定义的全连接层。初始化时，需要根据下游任务的需求，定义合适的输出维度。

6. 自定义模型训练

模型设计完成后，进入训练阶段。通过数据加载器 (DataLoader) 高效地批量处理数据，并使用优化器更新模型参数。

```
from torch.utils.data import DataLoader
from transformers import AdamW

# 实例化 DataLoader
```

```

data_loader = DataLoader(dataset, batch_size=16, shuffle=True)

# 初始化模型和优化器
model = SentimentAnalysisModel()
optimizer = AdamW(model.parameters(), lr=5e-5)

# 训练循环
for epoch in range(3): # 假设训练 3 个 epoch
    model.train()
    for batch in data_loader:
        optimizer.zero_grad()
        outputs = model(input_ids=batch['input_ids'],
            attention_mask=batch['attention_mask'])
        loss = nn.CrossEntropyLoss()(outputs, batch['labels'])
        loss.backward()
        optimizer.step()

```

6.1 数据加载

使用 `DataLoader` 实现批量数据加载。`DataLoader` 自动处理数据的批处理和随机打乱，确保训练的高效性和数据的多样性。

6.2 优化器

`AdamW` 是一种适用于 BERT 模型的优化器，结合了 Adam 和权重衰减的特点，能够有效地防止过拟合。

6.3 训练循环

训练循环包含前向传播（forward pass）、损失计算（loss calculation）、反向传播（backward pass）、参数更新（parameter update）等步骤。每个 epoch 都会对整个数据集进行一次遍历，更新模型参数。通常训练过程中会跟踪损失值的变化，以判断模型的收敛情况。

7. 最终效果评估与测试

在模型训练完成后，需要评估其在测试集上的性能。通常使用准确率、精确率、召回率和 F1 分数等指标来衡量模型的效果。

7.1 准确率 (Accuracy)

准确率是衡量分类模型整体性能的基本指标，计算公式为正确分类的样本数量除以总样本数量。

```

from sklearn.metrics import accuracy_score

# 假设有一个测试集 `test_dataset`，并且已经经过与训练集相同的预处理
test_loader = DataLoader(test_dataset, batch_size=16, shuffle=False)

# 评估模型在测试集上的性能
model.eval()
predictions, true_labels = [], []

with torch.no_grad():
    for batch in test_loader:
        outputs = model(input_ids=batch['input_ids'],
            attention_mask=batch['attention_mask'])
        _, preds = torch.max(outputs, dim=1)
        predictions.extend(preds.cpu().numpy())
        true_labels.extend(batch['labels'].cpu().numpy())

```

```
accuracy = accuracy_score(true_labels, predictions)
print(f"Accuracy: {accuracy:.4f}")
```

7.2 精确率、召回率和 F1 分数

精确率 (Precision) 和召回率 (Recall) 是分类模型的另两个重要指标，分别反映模型在正例预测上的精确性和召回能力。F1 分数是精确率和召回率的调和平均数，通常用于不平衡数据集的评估。

```
from sklearn.metrics import precision_score, recall_score, f1_score

precision = precision_score(true_labels, predictions, average='weighted')
recall = recall_score(true_labels, predictions, average='weighted')
f1 = f1_score(true_labels, predictions, average='weighted')

print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 Score: {f1:.4f}")
```

7.3 结果分析与模型优化

通过分析测试集上的结果，可以发现模型的强项和弱项。例如，如果 F1 分数较低，可能是由于数据集不平衡，导致模型在某些类别上表现不佳。通过调整超参数、改进数据预处理步骤，或使用更复杂的模型结构，可以进一步提高模型性能。

7.4 保存与加载模型

为了在未来使用训练好的模型，可以将其保存为文件，之后再加载进行推理或进一步的微调。

```
# 保存模型
torch.save(model.state_dict(), 'sentiment_analysis_model.pth')

# 加载模型
model = SentimentAnalysisModel()
model.load_state_dict(torch.load('sentiment_analysis_model.pth'))
model.eval()
```

8. 课件小结

在本课程中，我们详细介绍了如何使用 Hugging Face 的 BERT 模型进行中文评价情感分析的微调训练。我们从加载数据集、制作 Dataset、词汇表操作、模型设计、自定义训练，到最后的评估与测试，逐步讲解了整个微调过程。通过本课程，你需要掌握使用预训练语言模型进行下游任务微调的基本流程，并能应用到实际的 NLP 项目中。