

山东大学



实验 5: BERT 实践实验

大数据分析实践实验报告

姓	名:	郑坤武
学	号:	202200150184
班	级:	22 级公信班
学	院:	政治学与公共管理学院

2025 年 11 月 18 日

1 实验目的

本实验基于开源预训练语言模型 **BERT** 与 HuggingFace 生态 (**transformers** 与 **datasets**), 在 **GLUE-MRPC** 句子对复述识别任务上进行微调, 构建一个完整的“预训练模型 → 下游任务”教学范例。实验重点包括:

- 掌握 **BertTokenizerFast** 与 **BertForSequenceClassification** 的基本用法 (模型加载、分词编码、前向计算与损失获取)。
- 通过 **datasets** 库加载 GLUE 中的 **MRPC** 数据集, 完成句子对的批量编码与划分 (训练集 / 验证集 / 测试集)。
- 实现完整的 BERT 微调流程: **动态 padding 的 DataLoader**、训练循环、验证与测试评估 (Accuracy 与 F1), 理解预训练模型在下游任务中的迁移能力。
- 将微调后的模型与分词器保存到本地文件夹, 形成可复用的 **bert_mrpc_finetuned** 模型资源, 便于后续加载与推理。

2 实验环境

- 操作系统: macOS 12.x
- 开发工具: VS Code + 终端 (Terminal)
- 语言与运行环境: Python 3.12, 虚拟环境 **venv**
- 主要依赖库:
 - PyTorch (**torch**)
 - Transformers (**transformers**)
 - Datasets (**datasets**)
 - Scikit-learn (**scikit-learn**)
 - tqdm (**tqdm**)

3 具体实验步骤与结果分析

3.1 项目初始化与环境配置

在本地 macOS 上新建实验目录并创建 Python 虚拟环境, 使用 pip 安装所需依赖。终端中主要命令如下:

```
1 // 创建实验目录并进入
2 mkdir -p ~/bigdata_exp5_bert
3 cd ~/bigdata_exp5_bert
4
5 // 创建并激活虚拟环境
6 python3 -m venv venv
```

```
7 source venv/bin/activate
8
9 // 升级 pip 并安装依赖
10 pip install --upgrade pip
11 pip install torch transformers datasets scikit-learn tqdm
```

完成后, 在 VS Code 中打开该目录, 后续所有 Python 脚本 (如 `exp5_mrpc_bert.py`) 均放在此文件夹中, 方便统一管理。

3.2 数据集与任务介绍

本实验采用 GLUE (General Language Understanding Evaluation) 基准中的 **MRPC (Microsoft Research Paraphrase Corpus)** 数据集, 其核心是“句子对是否为语义复述”的二分类任务。MRPC 中每条样本由三个字段构成:

- `sentence1`: 第一个英文句子。
- `sentence2`: 第二个英文句子。
- `label`: 标签, 1 表示两句为语义等价 (复述), 0 表示不等价。

实验目标是: 给定一对句子 (s_1, s_2), 判断它们是否表达了相同的含义, 即构建一个基于 BERT 的句子对分类模型, 并在 MRPC 的验证集与测试集上评估其准确率与 F1 指标。

3.3 模型与分词器加载

在脚本中首先设置随机种子与设备选择逻辑, 保证实验可复现并优先利用硬件加速 (CUDA / MPS):

```
1 // 随机种子与设备选择
2 def set_seed(seed: int = 42):
3     random.seed(seed)
4     np.random.seed(seed)
5     torch.manual_seed(seed)
6     torch.cuda.manual_seed_all(seed)
7
8 def get_device() -> torch.device:
9     if torch.cuda.is_available():
10         print("使用设备: CUDA")
11         return torch.device("cuda")
12     if hasattr(torch.backends, "mps") and torch.backends.mps.is_available():
13         print("使用设备: MPS (Apple 芯片)")
14         return torch.device("mps")
15     print("使用设备: CPU")
16     return torch.device("cpu")
```

随后加载预训练好的 **BERT Base** 模型与对应分词器, 并移动到指定设备:

```
1 print("正在加载 GLUE MRPC 数据集...")
2 raw_datasets = load_dataset("glue", "mrpc")
```

```
3
4 print("正在加载 bert-base-uncased 模型...")
5 tokenizer = BertTokenizerFast.from_pretrained("bert-base-uncased")
6 model = BertForSequenceClassification.from_pretrained(
7     "bert-base-uncased", num_labels=2
8 )
9 model.to(device)
```

其中 num_labels=2 对应 MRPC 的二分类设置。

3.4 数据预处理与编码

由于 MRPC 是句子对任务，需要对两句输入联合编码。使用 BertTokenizerFast 提供的批处理接口对数据进行分词与截断：

```
1 def preprocess_function(examples, tokenizer, max_length: int = 128):
2     return tokenizer(
3         examples["sentence1"],
4         examples["sentence2"],
5         truncation=True,
6         max_length=max_length,
7     )
8
9 tokenized_datasets = raw_datasets.map(
10     lambda examples: preprocess_function(examples, tokenizer),
11     batched=True,
12     desc="Tokenization",
13 )
14
15 // 删去无关列，并重命名标签
16 tokenized_datasets = tokenized_datasets.remove_columns(
17     ["idx", "sentence1", "sentence2"]
18 )
19 tokenized_datasets = tokenized_datasets.rename_column("label", "labels")
```

这样模型在前向计算时即可直接接收 input_ids、attention_mask、token_type_ids 与 labels 等字段。

3.5 DataLoader 设计与动态 padding

为了兼顾效率与易用性，本实验采用 **动态 padding** 方式：不同 batch 的最大长度可以不同，同一 batch 内对齐到该批次的最大序列长度。

```
1 @dataclass
2 class Batch:
3     input_ids: torch.Tensor
4     attention_mask: torch.Tensor
5     token_type_ids: torch.Tensor
6     labels: torch.Tensor
7
```

```
8 def collate_fn(features: List[Dict], pad_token_id: int) -> Batch:
9     input_ids = [f["input_ids"] for f in features]
10    attention_mask = [f["attention_mask"] for f in features]
11    token_type_ids = [f["token_type_ids"] for f in features]
12    labels = torch.tensor([f["labels"] for f in features], dtype=torch. long)
13
14    max_len = max( len(x) for x in input_ids)
15
16    def pad(seqs, pad_value):
17        return torch.stack(
18            [
19                torch.tensor(
20                    list(seq) + [pad_value] * (max_len - len(seq)), dtype=torch. long
21                )
22                for seq in seqs
23            ]
24        )
25
26    input_ids = pad(input_ids, pad_token_id)
27    attention_mask = pad(attention_mask, 0)
28    token_type_ids = pad(token_type_ids, 0)
29
30    return Batch(
31        input_ids=input_ids,
32        attention_mask=attention_mask,
33        token_type_ids=token_type_ids,
34        labels=labels,
35    )
36
37 def create_dataloaders(tokenized_datasets, tokenizer, batch_size: int = 16):
38     pad_token_id = tokenizer.pad_token_id
39     train_loader = DataLoader(
40         tokenized_datasets["train"], batch_size=batch_size,
41         shuffle=True,
42         collate_fn=lambda batch: collate_fn(batch, pad_token_id),
43     )
44     eval_loader = DataLoader(
45         tokenized_datasets["validation"], batch_size=batch_size,
46         shuffle=False,
47         collate_fn=lambda batch: collate_fn(batch, pad_token_id),
48     )
49     test_loader = DataLoader(
50         tokenized_datasets["test"], batch_size=batch_size,
51         shuffle=False,
52         collate_fn=lambda batch: collate_fn(batch, pad_token_id),
53     )
54     return train_loader, eval_loader, test_loader
```

动态 padding 能有效减少无效的 [PAD] 计算, 提升训练速度和内存利用率。

3.6 训练策略与循环实现

训练阶段使用 AdamW 优化器与线性 warmup-decay 学习率调度器，训练轮数为 3 轮。关键超参数如下：

- 学习率： $2e-5$
- Batch size：16
- 训练轮数：3
- 优化器：torch.optim.AdamW
- 学习率调度：get_linear_schedule_with_warmup
- 梯度裁剪：max_norm=1.0

训练函数与评估函数的核心结构如下：

```
1 def train_one_epoch(model, dataloader, optimizer, scheduler, device, epoch_idx: int):
2     model.train()
3     running_loss = 0.0
4     progress_bar = tqdm(dataloader, desc=f"Epoch {epoch_idx+1} 训练中")
5
6     for batch in progress_bar:
7         batch = Batch(
8             input_ids=batch.input_ids.to(device),
9             attention_mask=batch.attention_mask.to(device),
10            token_type_ids=batch.token_type_ids.to(device),
11            labels=batch.labels.to(device),
12        )
13        outputs = model(
14            input_ids=batch.input_ids,
15            attention_mask=batch.attention_mask,
16            token_type_ids=batch.token_type_ids,
17            labels=batch.labels,
18        )
19        loss = outputs.loss
20        optimizer.zero_grad()
21        loss.backward()
22        torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
23        optimizer.step()
24        if scheduler is not None:
25            scheduler.step()
26        running_loss += loss.item()
27        progress_bar.set_postfix({"loss": f"{loss.item():.4f}"})
28    return running_loss / len(dataloader)
29
30 def evaluate(model, dataloader, device, desc: str = "验证集评估"):
31     model.eval()
32     all_labels, all_preds = [], []
33     with torch.no_grad():
```

```
34     progress_bar = tqdm(dataloader, desc=desc)
35     for batch in progress_bar:
36         batch = Batch(
37             input_ids=batch.input_ids.to(device),
38             attention_mask=batch.attention_mask.to(device),
39             token_type_ids=batch.token_type_ids.to(device),
40             labels=batch.labels.to(device),
41         )
42         outputs = model(
43             input_ids=batch.input_ids,
44             attention_mask=batch.attention_mask,
45             token_type_ids=batch.token_type_ids,
46         )
47         logits = outputs.logits
48         preds = torch.argmax(logits, dim=-1)
49         all_labels.extend(batch.labels.cpu().numpy().tolist())
50         all_preds.extend(preds.cpu().numpy().tolist())
51     acc = accuracy_score(all_labels, all_preds)
52     f1 = f1_score(all_labels, all_preds)
53     return acc, f1
```

每轮训练结束后在验证集上计算 Accuracy 与 F1，记录最佳 F1 对应的模型参数，用于最终测试集评估。

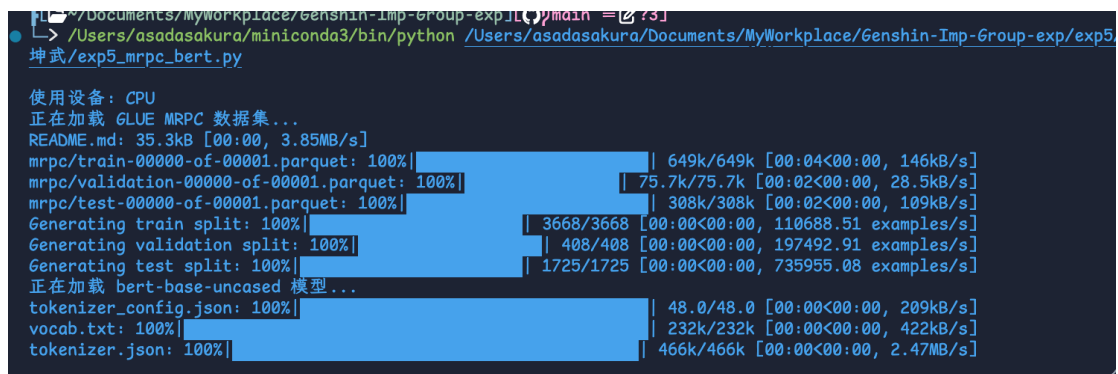
3.7 测试集评估与结果分析

在完成 3 轮训练后，将最佳验证效果对应的模型参数加载回来，在 MRPC 测试集上进行最终评估，得到如下指标：

- 测试集准确率： **0.8435**
- 测试集 F1 分数： **0.8854**

从结果可以看出，在没有特别复杂的调参与模型改造的情况下，单纯对 bert-base-uncased 进行少量 epoch 的微调，就可以在 MRPC 任务上取得较高的 F1 分数，说明 BERT 预训练模型已经在语义表征上学习到了丰富的语言知识，对小数据集任务具有良好的迁移能力。

如图 2 所示，展示训练过程与最终测试结果。



```
使用设备: CPU
正在加载 GLUE MRPC 数据集...
README.md: 35.3kB [00:00, 3.85MB/s]
mrpc/train-00000-of-00001.parquet: 100%|          | 649k/649k [00:04<00:00, 146kB/s]
mrpc/validation-00000-of-00001.parquet: 100%|          | 75.7k/75.7k [00:02<00:00, 28.5kB/s]
mrpc/test-00000-of-00001.parquet: 100%|          | 308k/308k [00:02<00:00, 109kB/s]
Generating train split: 100%|          | 3668/3668 [00:00<00:00, 110688.51 examples/s]
Generating validation split: 100%|          | 408/408 [00:00<00:00, 197492.91 examples/s]
Generating test split: 100%|          | 1725/1725 [00:00<00:00, 735955.08 examples/s]
正在加载 bert-base-uncased 模型...
tokenizer_config.json: 100%|          | 48.0/48.0 [00:00<00:00, 209kB/s]
vocab.txt: 100%|          | 232k/232k [00:00<00:00, 422kB/s]
tokenizer.json: 100%|          | 466k/466k [00:00<00:00, 2.47MB/s]
```

图 1: 终端中 BERT 微调过程示例

```
~/Documents/MyWorkplace/Genshin-Imp-Group-exp/finetune [2] 13
> /Users/asadasakura/miniconda3/bin/python /Users/asadasakura/Documents/MyWorkplace/Genshin-Imp-Group-exp/exp5_
坤武/exp5_mrpc_bert.py
Epoch 2/3 平均训练损失: 0.3176
验证集评估: 100% | 26/26 [00:09<00:00, 2.70it/s]
验证集准确率: 0.8260, 验证集 F1 分数: 0.8632
Epoch 3 训练中: 100% | 230/230 [06:02<00:00, 1.58s/it, loss=0.0101]
Epoch 3/3 平均训练损失: 0.1567
验证集评估: 100% | 26/26 [00:09<00:00, 2.67it/s]
验证集准确率: 0.8603, 验证集 F1 分数: 0.9005
测试集评估: 100% | 108/108 [00:40<00:00, 2.68it/s]
在测试集上的最终结果:
测试集准确率: 0.8435
测试集 F1 分数: 0.8854
微调后的模型已保存到: /Users/asadasakura/Documents/MyWorkplace/Genshin-Imp-Group-exp/bert_mrpc_finetuned
你可以在本地文件管理器中打开该文件夹查看。用于实验报告
```

图 2: 终端中 BERT 微调测试集评估结果示例

3.8 模型保存与本地文件结构

实验最后将微调后的模型与分词器保存到本地目录 `bert_mrpc_finetuned` 中:

```
1 save_dir = "bert_mrpc_finetuned"
2 os.makedirs(save_dir, exist_ok=True)
3 model.save_pretrained(save_dir)
4 tokenizer.save_pretrained(save_dir)
5 print(f"微调后的模型已保存到: {os.path.abspath(save_dir)}")
```

在本地实际运行中, 模型保存路径为:

`/Users/.../bert_mrpc_finetuned`

该目录中包含配置文件、权重文件与分词器相关文件等, 可通过本地图形界面截图展示, 如图 3 所示。

图 3: 本地 `bert_mrpc_finetuned` 目录结构示意

通过该步骤, 验证了微调模型可以顺利落地到本地文件系统, 为后续加载、推理或继续微调提供了基础。

4 实验总结与收获

本实验从零开始在本地 macOS 环境中,基于 HuggingFace 生态完成了 BERT 在 GLUE-MRPC 任务上的微调,实现了“预训练模型 + 下游任务”的典型范例。整个过程涵盖了数据加载与预处理、动态 padding 的 DataLoader 设计、训练与验证评估、测试集性能分析以及模型持久化保存等关键环节。

通过本次实验,我对以下几点有了更深刻的理解:

- **预训练模型的优势:** 与从头训练模型相比,直接在 `bert-base-uncased` 上做少量微调,就能在小样本任务上取得较高的准确率与 F1,体现出大规模预训练在表示学习上的强大能力。
- **工程化流程意识:** 通过自定义 `collate_fn`、设计训练 / 评估函数与保存模型等步骤,初步建立起一个可复用的 NLP 训练脚手架,后续可以方便迁移到其他文本分类、句子匹配等任务上。
- **调参与扩展空间:** 当前实验仅使用了默认的 BERT Base 与简单的超参数配置,后续可以尝试更深层的模型(如 RoBERTa、DeBERTa 等)、不同学习率 / batch size、更多 epoch 或加入早停(early stopping)机制,以探索性能进一步提升的可能性。

总体而言,本实验帮助我从“会用预训练模型”迈向“能够自己搭建完整微调流程”的阶段,为后续学习更复杂的 NLP 任务(如问答、文本生成、多任务学习等)打下了实践基础。