# NLP practice pj：基于预训练transformer模型的wnli分类

本次pj是 基于预训练transformer模型的wnli分类，我最终通过更换模型为 `microsoft/deberta-v3-large` ,使用warmup和cosine策略，以及调整lr，bs，epoch等超参，实现了test的score达到85.6的效果
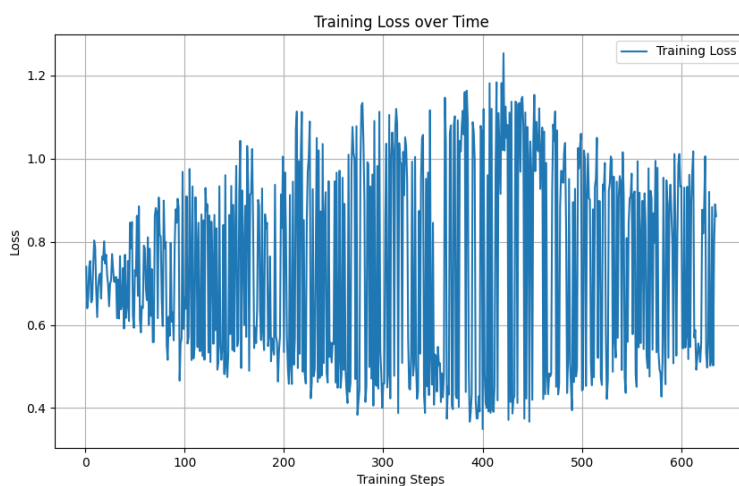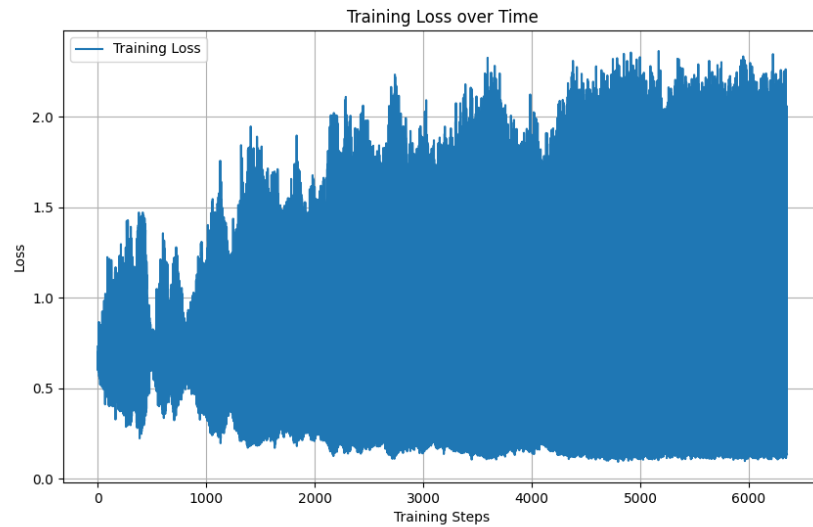
具体的PJ实验内容如下

首先我使用了原模型prajjwal1/bert-tiny进行实验

原模型训练之后的结果：

```
1  'eval_loss': 0.8087890148162842, 'eval_accuracy': 0.43661971830985913
```
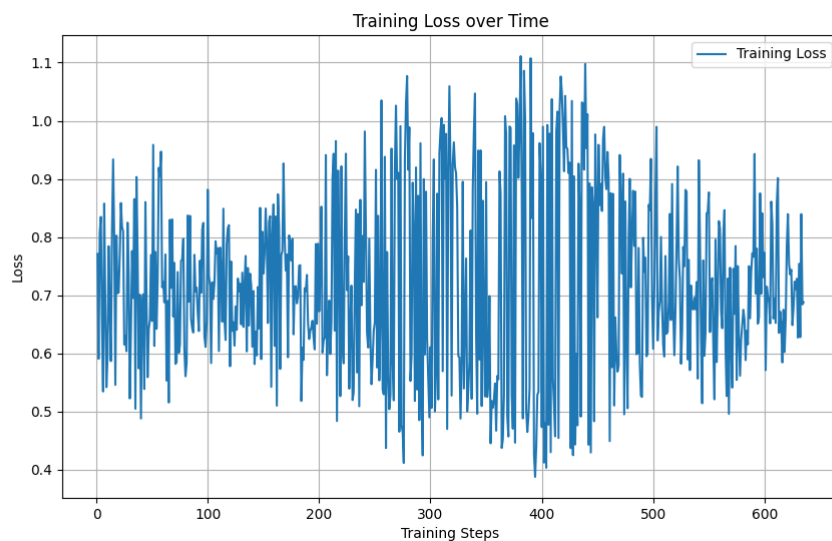
loss曲线



如果设置epoch为10，得到的训练曲线会这样

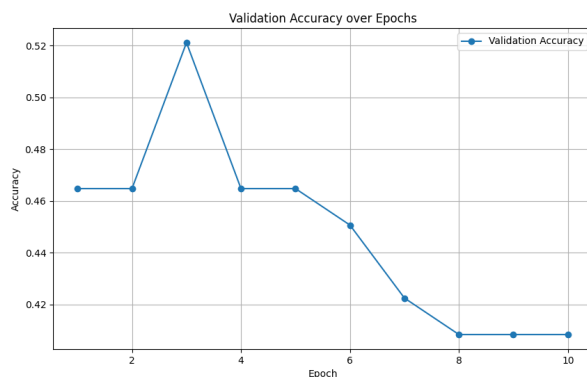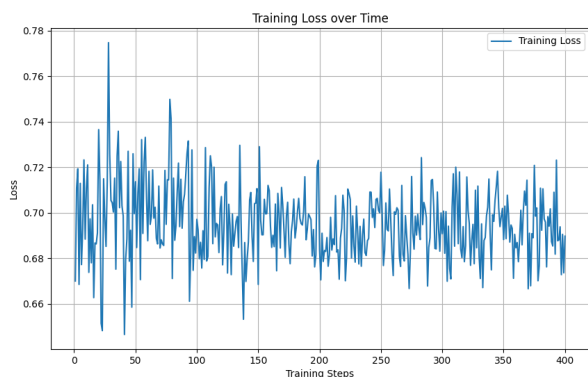Training Loss over Time

由曲线看出，模型并没有收敛，可能是学习率设置太大或者是学习率调度器有问题

改进

- 增加了warm up预热策略 ，设置warmup_ratio=0.2，在训练开始的 20% 时间内逐渐提高学习率至设定值。该预热策略可以在训练初期防止梯度爆炸，提高训练的稳定性和最终效果

- 学习率调度器改为余弦退火lr_scheduler_type="cosine"，这种调度方法可以帮助模型更快收敛，并在后期进一步稳定



Training Loss over Time

对比之前在epoch快结束的时候可以观测到loss曲线有收敛的趋势

将epoch设置为10，得到的曲线如下

loss曲线虽然收敛，但是没有明显的下降趋势，valid的accuracy在epoch=3的时候达到峰值，之后逐渐降低，应该是出现了过拟合的现象，最高值是accuracy=0.52

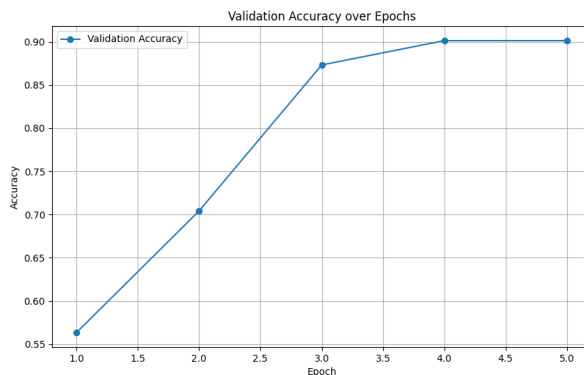尝试修改学习率和bs，也没有得到很好的效果，最高的accuracy只有0.573

因此我觉得可能的原因在于数据和模型，主要是亮点：

- 数据量太少或太局部以至于模型训练效果不好，无法泛化

- 模型本身性能不好

数据集的扩大较为困难，因此我尝试将模型更换

表格记录了使用的尝试过的模型以及在验证集的最好的效果

| model | Best valid accuracy | bs | lr | epoch |
|---|---|---|---|---|
| prajjwal1/bert-tiny | 0.573 | 16 | 2e-5 | 2 |
| microsoft/deberta-v3-small | 0.56 | 16 | 2e-5 | 3 |
| albert/albert-base-v2 | 0.56 | 16 | 2e-5 | 4 |
| xlnet/xlnet-base-cased | 0.56 | 16 | 2e-5 | 3 |
| microsoft/deberta-v3-large | 0.91 | 16 | 2e-5 | 5 |

从表格中得出，使用 `microsoft/deberta-v3-large` 效果最好，训练5个epoch的loss曲线和验证集的accuracy如下：

测试结果：test accuracy score = 85.6

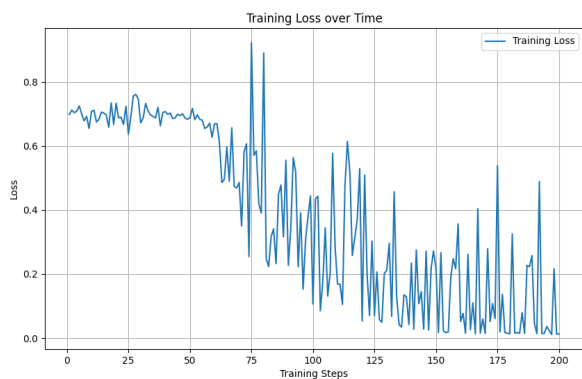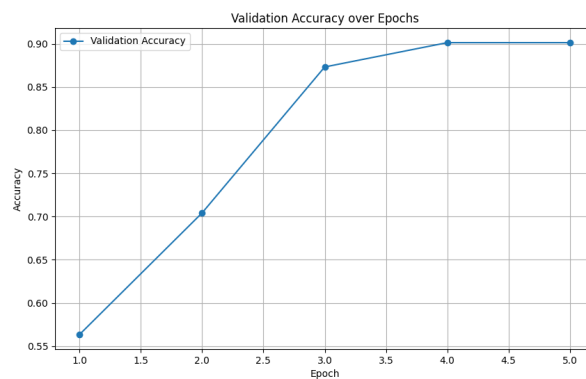| Score: 58.8 | | |
|---|---|---|
| **PRIMARY** | **DIAGNOSTICS** | |
| **Task** | **Metric** | **Score** |
| The Corpus of Linguistic Acceptability | Matthew's Corr | 0.0 |
| The Stanford Sentiment Treebank | Accuracy | 80.0 |
| Microsoft Research Paraphrase Corpus | F1 / Accuracy | 81.5/73.4 |
| Semantic Textual Similarity Benchmark | Pearson-Spearman Corr | 61.2/59.1 |
| Quora Question Pairs | F1 / Accuracy | 51.4/79.1 |
| MultiNLI Matched | Accuracy | 56.0 |
| MultiNLI Mismatched | Accuracy | 56.4 |
| Question NLI | Accuracy | 50.4 |
| Recognizing Textual Entailment | Accuracy | 54.1 |
| Winograd NLI | Accuracy | 85.6 |
| Diagnostics Main | Matthew's Corr | 9.2 |

测试结果的链接：

https://gluebenchmark.com/submission/uW4I3J9YPZRIaLXWwlipxvcl9dh1/-OACoSSW1POtzrjU_GMH

最后，我做了一组消融实验来比对学习率调度方式对模型性能的影响
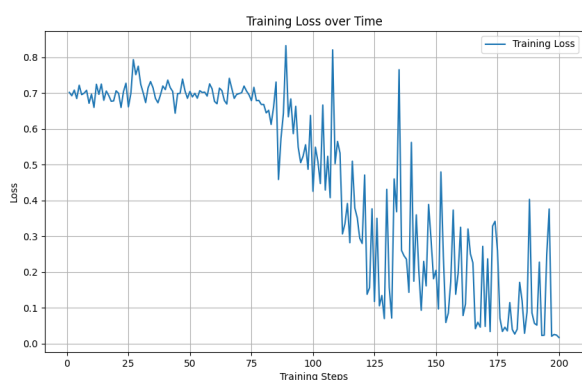
以下是实验结果，由上至下分别为，使用了cosine+warmup策略，使用了linear+warmup策略，只使用linear策略，这三种方式的模型训练曲线图，左图为loss曲线，右图为每个epoch的valid accuracy曲线
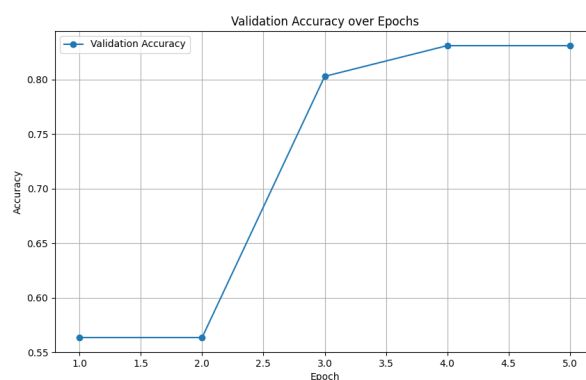
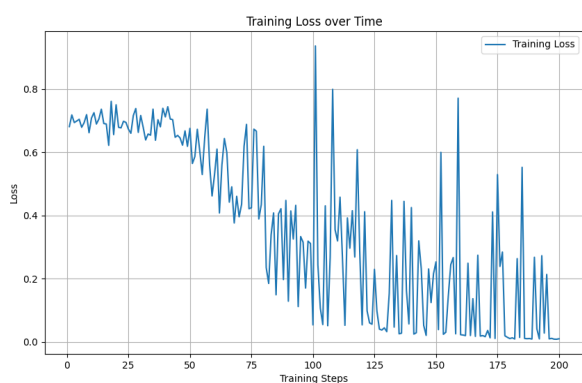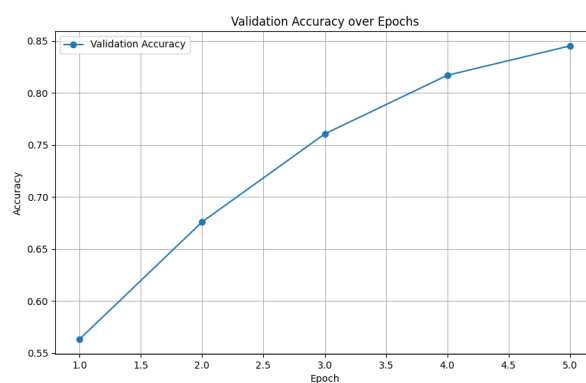cosine+warm up



cosine+warmup



linear+warmup



linear+warmup



linear



linear

从上述图中可以看出，使用了warmup策略后模型训练可以更快的收敛，在第3个epoch的时候就已经达到了0.8以上的accuracy，使用了cosine策略之后模型的性能进一步提升，accuracy达到了0.9以上

总结

本次实验，我使用了warmup和cosine策略，使用了更换模型的方法，顺利完成了本次pj

以下贴一个config.json

```
1  {
```

```
 2      "_name_or_path": "microsoft/deberta-v3-large",
 3      "architectures": [
 4        "DebertaV2ForSequenceClassification"
 5      ],
 6      "attention_probs_dropout_prob": 0.1,
 7      "hidden_act": "gelu",
 8      "hidden_dropout_prob": 0.1,
 9      "hidden_size": 1024,
10      "initializer_range": 0.02,
11      "intermediate_size": 4096,
12      "layer_norm_eps": 1e-07,
13      "max_position_embeddings": 512,
14      "max_relative_positions": -1,
15      "model_type": "deberta-v2",
16      "norm_rel_ebd": "layer_norm",
17      "num_attention_heads": 16,
18      "num_hidden_layers": 24,
19      "pad_token_id": 0,
20      "pooler_dropout": 0,
21      "pooler_hidden_act": "gelu",
22      "pooler_hidden_size": 1024,
23      "pos_att_type": [
24        "p2c",
25        "c2p"
26      ],
27      "position_biased_input": false,
28      "position_buckets": 256,
29      "relative_attention": true,
30      "share_att_key": true,
31      "torch_dtype": "float32",
32      "transformers_version": "4.45.2",
33      "type_vocab_size": 0,
34      "vocab_size": 128100
35  }
36
```

训练超参

- lr=2e-5

- epoch=5

- bs=16

- lr_scheduler_type="cosine"

- warmup_ratio=0.2

- gradient_accumulation_steps=1

test的结果链接

https://gluebenchmark.com/submission/uW4I3J9YPZRIaLXWwlipxvcl9dh1/-OACoSSW1POtzrjU_GMH

train使用的代码

```python
import argparse

import evaluate
import numpy as np
from datasets import load_dataset
from transformers import (
    AutoModelForSequenceClassification,
    AutoTokenizer,
    EvalPrediction,
    Trainer,
    TrainingArguments,
)
import matplotlib.pyplot as plt

def main(args):
    dataset = load_dataset("SetFit/wnli")
    tokenizer = AutoTokenizer.from_pretrained(args.model,
    trust_remote_code=True)

    def preprocess_function(examples):
        return {
            **tokenizer(examples["text1"], examples["text2"]),
            "label": examples["label"],
        }

    dataset = dataset.map(preprocess_function, batched=True)
    print(tokenizer.decode(dataset["train"]["input_ids"][0]))

    model = AutoModelForSequenceClassification.from_pretrained(
        args.model,
        num_labels=2,
        trust_remote_code=True,
    )
    print(sum(p.numel() for p in model.parameters()))

    metric = evaluate.load("accuracy")

    def compute_metrics(p: EvalPrediction):
        preds = np.argmax(p.predictions, axis=-1)
        result = metric.compute(predictions=preds, references=p.label_ids)
```

```
40          return result
41
42      trainer = Trainer(
43          model=model,
44          tokenizer=tokenizer,
45          train_dataset=dataset["train"],
46          eval_dataset=dataset["validation"],
47          compute_metrics=compute_metrics,
48          args=TrainingArguments(
49              output_dir=args.output,
50              eval_strategy="epoch",
51              save_strategy="epoch",
52              logging_steps=1,
53              learning_rate=args.lr,
54              per_device_train_batch_size=args.bs,
55              per_device_eval_batch_size=args.bs,
56              gradient_accumulation_steps=args.accum,
57              num_train_epochs=args.epoch,
58              lr_scheduler_type="cosine",
59              warmup_ratio=0.2,
60          ),
61      )
62      trainer.train()
63
64      # 提取损失值和准确率
65      losses = [log["loss"] for log in trainer.state.log_history if "loss" in
   log]
66      accuracies = [log["eval_accuracy"] for log in trainer.state.log_history if
   "eval_accuracy" in log]
67      epochs = [log["epoch"] for log in trainer.state.log_history if
   "eval_accuracy" in log]
68
69      # 绘制损失曲线并保存
70      plt.figure(figsize=(10, 6))
71      plt.plot(range(1, len(losses) + 1), losses, label="Training Loss")
72      plt.xlabel("Training Steps")
73      plt.ylabel("Loss")
74      plt.title("Training Loss over Time")
75      plt.legend()
76      plt.grid()
77      plt.savefig("training_loss_curve.png")
78
79      # 绘制准确率曲线并保存
80      plt.figure(figsize=(10, 6))
81      plt.plot(epochs, accuracies, label="Validation Accuracy", marker='o')
82      plt.xlabel("Epoch")
83      plt.ylabel("Accuracy")
```

```python
84      plt.title("Validation Accuracy over Epochs")
85      plt.legend()
86      plt.grid()
87      plt.savefig("validation_accuracy_curve.png")
88
89  if __name__ == "__main__":
90      parser = argparse.ArgumentParser()
91      parser.add_argument("--model", type=str, default="prajjwal1/bert-tiny")
92      parser.add_argument("--lr", type=float, default=2e-5)
93      parser.add_argument("--epoch", type=int, default=1)
94      parser.add_argument("--bs", type=int, default=1)
95      parser.add_argument("--accum", type=int, default=1)
96      parser.add_argument("--output", type=str, default="output")
97      args = parser.parse_args()
98      main(args)
99
```