Date: 2023.12.10

Notes: The PyTorch template may be updated in the future, but the code in this example may not be affected by those changes.

# SST2

This is the record of how to tweak the pytorch template for this project, as well as what the training procedure looks like.

## 1. Tweak code

Remove the code unneeded, e.g., the code for cv.

### 1.1 About data

The SST2 dataset will be loaded from huggingface. Specify the dataset in **./main.py** and load the train, valid, and test splits.

```
48      # load data from huggingface
49      cache_dir = "././huggingface"
50   💡 dataset_path = "SetFit/sst2"
51      raw_dataset = load_dataset(path=dataset_path, cache_dir=cache_dir)
52
53      train_data = raw_dataset['train']
54      valid_data = raw_dataset['validation']
55      test_data = raw_dataset['test']
56
```

The default function for preprocessing in **./preprocess.py** is fine, so I didn't touch it.

Also, the TextDataset class in **./dataset.py** could be used directly.

### 1.2 About model

The default MyModel in **./model.py** is designed for text classification task, I invoked it without any changes.

### 1.3 About training

By default, the template uses CrossEntropyLoss for criterion, AdamW for optimizer, CosineAnnealingWarmRestarts for lr scheduler, which seems appropriate. So I didn't touch these in **./main.py**, either.

The Trainer in **./trainer.py** is ready-to-use, and it is recommended to use it directly without any alterations.

The template includes accuracy for test method, which just fit my demand in this simple project. So I kept it and didn't add more test methods.

## 1.4 About config

Tweak configurations in **./config.yaml**:

Use wandb to track experiment, set related config.

```yaml
1   seed: 6
2   use_wandb: True
3
4   # config for wandb
5   wandb_cfg:
6     project: "SST2"
7     notes: "training details on the process of global rank 0"
8     tags: ["SST2", "TransformerEncoder"]
9     watch_model: True
10    # required if `watch_model` is True
11    watch_model_freq: 1
12
```

Tweak config for preprocess, here the hyperparams are about tokenizer and vocabulary.

```yaml
13    # config for NLP preprocess
14    preprocess_nlp_cfg:
15      lowercase: True
16      rm_punctuation: True
17      rm_stopword: False
18      lemmatization: True
19      min_freq: 3
20      max_tokens: 10000
21
```

Tweak config for dataloader (e.g., batch_size), model (in this case, is a TransformerEncoder), optimizer (e.g., lr), and lr scheduler.
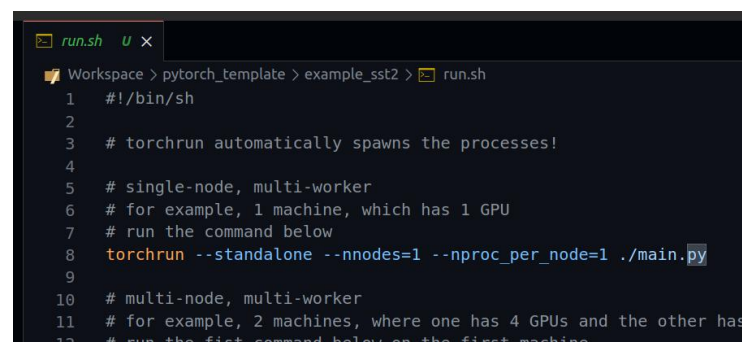
```yaml
22    # config for loader
23    loader_cfg:
24      batch_size: 32
25      num_workers: 24
26      pin_memory: True
27
28    # config for model
29    model_cfg:
30      vocab_size: 10000
31      embed_dim: 128
32      nhead: 2
33      dim_feedforward: 512
34      num_layers: 1
35      num_classes: 2
36      dropout: 0.1
37
38    # config for optimizer
39    optimizer_cfg:
40      lr: 0.001
41      weight_decay: 0.01
42
43    # config for scheduler
44    scheduler_cfg:
45      T_0: 4
46      T_mult: 2
47
```

Tweak config for training. Here I made it to:

train up to 10 epochs;
not use gradient accumulation;
do validation and test accuracy;
save logs, best model, and checkpoints during training;
train from scratch rather than from checkpoint;
start validation at epoch 1 and at every 1 epoch;
start testing accuracy at epoch 1 and at every 1 epoch;
save logs and checkpoints to an existed directoty;
use accuracy to measure best model;
save latest checkpoint and checkpoints at specified epochs.

```yaml
48   # config for train
49   train_cfg:
50     max_epoch: 10
51     accum_step: 1
52     do_valid: True
53     do_test: True
54     save_log: True
55     save_best: True
56     save_checkpoint: True
57     resume_checkpoint: False
58     # required if `do_valid` is True
59     valid_start: 1
60     valid_step: 1
61     # required if `do_test` is True
62     test_start: 1
63     test_step: 1
64     # required if `save_*` is True
65     save_dir: "./sst2_ckpt"
66     # required if `save_best` is True
67     measure_best: "accuracy"
68     measure_mode: "max"
69     # required if `save_checkpoint` is True
70     checkpoint_latest: True
71     checkpoint_list: [4, 8]
72     # required if `resume_checkpoint` is True
73     resume_path: null
```

Finally, adjust **./run.sh** based on the machine architecture. I ran it on my laptop with single gpu.

```sh
run.sh   U ×

Workspace > pytorch_template > example_sst2 > run.sh
1    #!/bin/sh
2
3    # torchrun automatically spawns the processes!
4
5    # single-node, multi-worker
6    # for example, 1 machine, which has 1 GPU
7    # run the command below
8    torchrun --standalone --nnodes=1 --nproc_per_node=1 ./main.py
9
10   # multi-node, multi-worker
11   # for example, 2 machines, where one has 4 GPUs and the other has
12   # run the fist command below on the first machine
```

# 2. Training appearance

terminal:

```
in scores: accuracy: 0.9017341040462428 | Valid scores: accuracy: 0.7488532110091743 | Time/epoch: 21.02294 seconds
2023-12-10 16:45:40 - INFO - New best model: valid accuracy update from 0.7293577981651376 to 0.7488532110091743
2023-12-10 16:45:40 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch9.pth ...
100%|                                                                    | 217/217 [00:09<00:00, 22.08it/s]
100%|                                                                    | 28/28 [00:01<00:00, 23.93it/s]
100%|                                                                    | 217/217 [00:09<00:00, 22.78it/s]
100%|                                                                    | 28/28 [00:01<00:00, 22.94it/s]
2023-12-10 16:46:01 - INFO - [GPU0] | Epoch 10/10 | Train loss: 0.4105628948607203 | Valid loss: 0.5536428679312978 | Tr
ain scores: accuracy: 0.9134393063583816 | Valid scores: accuracy: 0.7534403669724771 | Time/epoch: 21.74936 seconds
2023-12-10 16:46:01 - INFO - New best model: valid accuracy update from 0.7488532110091743 to 0.7534403669724771
2023-12-10 16:46:01 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch10.pth ...
2023-12-10 16:46:01 - INFO - ---------- End of training. Total time: 198.88231 seconds ----------
wandb:
wandb:
wandb: Run history:
wandb:               epoch
wandb: eval/best_valid_accuracy
wandb:     eval/train_accuracy
wandb:     eval/valid_accuracy
wandb:          eval/valid_loss
wandb:       train/epoch_time
wandb:               train/lr
wandb:       train/train_loss
wandb:
wandb: Run summary:
```



```
wandb:       train/epoch_time
wandb:               train/lr
wandb:       train/train_loss
wandb:
wandb: Run summary:
wandb:               epoch 10
wandb: eval/best_valid_accuracy 0.75344
wandb:     eval/train_accuracy 0.91344
wandb:     eval/valid_accuracy 0.75344
wandb:          eval/valid_loss 0.55364
wandb:       train/epoch_time 21.74936
wandb:               train/lr 0.00015
wandb:       train/train_loss 0.41056
wandb:
wandb: 🚀 View run clear-cloud-8 at: https://wandb.ai/nehc0/SST2/runs/q3wnn9v5
wandb: ⚡ View job at https://wandb.ai/nehc0/SST2/jobs/QXJ0aWZhY3RDb2xsZWN0aW9uOjEyMjMzODcwMw==/version_details/v4
wandb: Synced 6 W&B file(s), 0 media file(s), 0 artifact file(s) and 0 other file(s)
wandb: Find logs at: ./wandb/run-20231210_164239-q3wnn9v5/logs
2023-12-10 16:46:11 - INFO - Loading checkpoint: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch10.pth ...
2023-12-10 16:46:11 - INFO - Checkpoint loaded successfully.
100%|                                                                    | 57/57 [00:01<00:00, 46.28it/s]
2023-12-10 16:46:12 - INFO - Scores on test dataset: accuracy: 0.7358594179022515
~/Workspace/pytorch_template/example_sst2  on main +6 !27 ?4 ················ ✔ took 4m 17s  dl_pytorch 
```

wandb:

saved logs and checkpoints:

```
example_sst2
  .huggingface
sst2_ckpt
  run@231210_16:42:42
    best_model_epoch10.pth
    checkpoint_epoch4.pth
    checkpoint_epoch8.pth
    latest_checkpoint_epoch10.pth
    train.log
```



```
train.log U ✕
Workspace > pytorch_template > example_sst2 > sst2_ckpt > run@231210_16:42:42 > train.log
  1  2023-12-10 16:42:42 - INFO - ---------------- config ----------------
  2  seed: 6
  3  use_wandb: True
  4  wandb_cfg: {
  5  project: SST2
  6  notes: training details on the process of global rank 0
  7  tags: ['SST2', 'TransformerEncoder']
  8  watch_model: True
  9  watch_model_freq: 1
 10  }
 11  preprocess_nlp_cfg: {
 12  lowercase: True
 13  rm_punctuation: True
 14  rm_stopword: False
 15  lemmatization: True
 16  min_freq: 3
 17  max_tokens: 10000
 18  }
 19  loader_cfg: {
 20  batch_size: 32
 21  num_workers: 24
 22  pin_memory: True
 23  batch_size_per_proc: 32
 24  effective_batch_size: 32
 25  }
```



```
train.log U ✕
Workspace > pytorch_template > example_sst2 > sst2_ckpt > run@231210_16:42:42 > train.log
 56  save_dir: ./sst2_ckpt
 57  measure_best: accuracy
 58  measure_mode: max
 59  checkpoint_latest: True
 60  checkpoint_list: [4, 8]
 61  resume_path: None
 62  }
 63  world_size: 1
 64
 65  2023-12-10 16:42:42 - INFO - ---------- Start of training. Good day! ----------
 66  2023-12-10 16:43:01 - INFO - [GPU0] | Epoch 1/10 | Train loss: 0.6704142692451653 | Valid loss: 0.6239489380802427 | Train
 67  2023-12-10 16:43:01 - INFO - New best model: valid accuracy update from -inf to 0.6628440366972477
 68  2023-12-10 16:43:01 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch1.pth ...
 69  2023-12-10 16:43:21 - INFO - [GPU0] | Epoch 2/10 | Train loss: 0.6001296817981703 | Valid loss: 0.5849213749170303 | Train
 70  2023-12-10 16:43:21 - INFO - New best model: valid accuracy update from 0.6628440366972477 to 0.7144495412844036
 71  2023-12-10 16:43:21 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch2.pth ...
 72  2023-12-10 16:43:40 - INFO - [GPU0] | Epoch 3/10 | Train loss: 0.5372536731755129 | Valid loss: 0.5736100247928074 | Train
 73  2023-12-10 16:43:40 - INFO - New best model: valid accuracy update from 0.7144495412844036 to 0.7282110091743119
 74  2023-12-10 16:43:40 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch3.pth ...
 75  2023-12-10 16:43:59 - INFO - [GPU0] | Epoch 4/10 | Train loss: 0.4971857097017051 | Valid loss: 0.5704488051789147 | Train
 76  2023-12-10 16:43:59 - INFO - New best model: valid accuracy update from 0.7282110091743119 to 0.7293577981651376
 77  2023-12-10 16:43:59 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch4.pth ...
 78  2023-12-10 16:43:59 - INFO - Saving checkpoint: ./sst2_ckpt/run@231210_16:42:42/checkpoint_epoch4.pth ...
 79  2023-12-10 16:44:18 - INFO - [GPU0] | Epoch 5/10 | Train loss: 0.5088460817589738 | Valid loss: 0.5786607595426696 | Train
 80  2023-12-10 16:44:38 - INFO - [GPU0] | Epoch 6/10 | Train loss: 0.4883294870501839 | Valid loss: 0.5736998000315258 | Train
 81  2023-12-10 16:44:58 - INFO - [GPU0] | Epoch 7/10 | Train loss: 0.4580648419219777 | Valid loss: 0.5765066466161183 | Train
 82  2023-12-10 16:45:18 - INFO - [GPU0] | Epoch 8/10 | Train loss: 0.4451347106063421 | Valid loss: 0.5650713560836655 | Train
 83  2023-12-10 16:45:18 - INFO - Saving checkpoint: ./sst2_ckpt/run@231210_16:42:42/checkpoint_epoch8.pth ...
 84  2023-12-10 16:45:40 - INFO - [GPU0] | Epoch 9/10 | Train loss: 0.4242527932615324 | Valid loss: 0.5551174240452903 | Train
 85  2023-12-10 16:45:40 - INFO - New best model: valid accuracy update from 0.7293577981651376 to 0.7488532110091743
 86  2023-12-10 16:45:40 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch9.pth ...
 87  2023-12-10 16:46:01 - INFO - [GPU0] | Epoch 10/10 | Train loss: 0.4105628948607203 | Valid loss: 0.5536428679312978 | Trai
 88  2023-12-10 16:46:01 - INFO - New best model: valid accuracy update from 0.7488532110091743 to 0.7534403669724771
 89  2023-12-10 16:46:01 - INFO - Saving best model: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch10.pth ...
 90  2023-12-10 16:46:01 - INFO - ---------- End of training. Total time: 198.88231 seconds ----------
 91  2023-12-10 16:46:11 - INFO - Loading checkpoint: ./sst2_ckpt/run@231210_16:42:42/best_model_epoch10.pth ...
 92  2023-12-10 16:46:11 - INFO - Checkpoint loaded successfully.
 93  2023-12-10 16:46:12 - INFO - Scores on test dataset: accuracy: 0.7358594179022515
 94
```