

CodeBlame82 lines (72 loc) · 2.83 KB

Raw

```
1  from utils import *
2  from crosscoder import CrossCoder
3  from buffer import Buffer
4  import tqdm
5
6  from torch.nn.utils import clip_grad_norm_
7  class Trainer:
8      def __init__(self, cfg, model_A, model_B, all_tokens):
9          self.cfg = cfg
10         self.model_A = model_A
11         self.model_B = model_B
12         self.crosscoder = CrossCoder(cfg)
13         self.buffer = Buffer(cfg, model_A, model_B, all_tokens)
14         self.total_steps = cfg["num_tokens"] // cfg["batch_size"]
15
16         self.optimizer = torch.optim.Adam(
17             self.crosscoder.parameters(),
18             lr=cfg["lr"],
19             betas=(cfg["beta1"], cfg["beta2"]),
20         )
21         self.scheduler = torch.optim.lr_scheduler.LambdaLR(
22             self.optimizer, self.lr_lambda
23         )
24         self.step_counter = 0
25
```

Symbols

Find definitions and references for functions and other symbols in this file by clicking a symbol below or in the code.

Filter symbols

class Trainer

func __init__

func lr_lambda

func get_l1_coeff

func step

func log

func save

func train

```

26     wandb.init(project=cfg["wandb_project"], entity=cfg["wandb_entity"])
27
28     def lr_lambda(self, step):
29         if step < 0.8 * self.total_steps:
30             return 1.0
31         else:
32             return 1.0 - (step - 0.8 * self.total_steps) / (0.2 * self.total_steps)
33
34     def get_l1_coeff(self):
35         # Linearly increases from 0 to cfg["l1_coeff"] over the first 0.05 * self.t
36         if self.step_counter < 0.05 * self.total_steps:
37             return self.cfg["l1_coeff"] * self.step_counter / (0.05 * self.total_st
38         else:
39             return self.cfg["l1_coeff"]
40
41     def step(self):
42         acts = self.buffer.next()
43         losses = self.crosscoder.get_losses(acts)
44         loss = losses.l2_loss + self.get_l1_coeff() * losses.l1_loss
45         loss.backward()
46         clip_grad_norm_(self.crosscoder.parameters(), max_norm=1.0)
47         self.optimizer.step()
48         self.scheduler.step()
49         self.optimizer.zero_grad()
50
51         loss_dict = {
52             "loss": loss.item(),
53             "l2_loss": losses.l2_loss.item(),
54             "l1_loss": losses.l1_loss.item(),
55             "l0_loss": losses.l0_loss.item(),
56             "l1_coeff": self.get_l1_coeff(),
57             "lr": self.scheduler.get_last_lr()[0],
58             "explained_variance": losses.explained_variance.mean().item(),
59             "explained_variance_A": losses.explained_variance_A.mean().item(),
60             "explained_variance_B": losses.explained_variance_B.mean().item(),
61         }
62         self.step_counter += 1
63         return loss_dict
64

```

```
65     def log(self, loss_dict):
66         wandb.log(loss_dict, step=self.step_counter)
67         print(loss_dict)
68
69     def save(self):
70         self.crosscoder.save()
71
72     def train(self):
73         self.step_counter = 0
74         try:
75             for i in tqdm.trange(self.total_steps):
76                 loss_dict = self.step()
77                 if i % self.cfg["log_every"] == 0:
78                     self.log(loss_dict)
79                 if (i + 1) % self.cfg["save_every"] == 0:
80                     self.save()
81         finally:
82             self.save()
```