

CodeBlame115 lines (103 loc) · 4.68 KB

Raw

```
1  from utils import *
2  from transformer_lens import ActivationCache
3  import tqdm
4
5  class Buffer:
6      """
7      This defines a data buffer, to store a stack of acts across both model that car
8      """
9
10 def __init__(self, cfg, model_A, model_B, all_tokens):
11     assert model_A.cfg.d_model == model_B.cfg.d_model
12     self.cfg = cfg
13     self.buffer_size = cfg["batch_size"] * cfg["buffer_mult"]
14     self.buffer_batches = self.buffer_size // (cfg["seq_len"] - 1)
15     self.buffer_size = self.buffer_batches * (cfg["seq_len"] - 1)
16     self.buffer = torch.zeros(
17         (self.buffer_size, 2, model_A.cfg.d_model),
18         dtype=torch.bfloat16,
19         requires_grad=False,
20     ).to(cfg["device"]) # hardcoding 2 for model diffing
21     self.cfg = cfg
22     self.model_A = model_A
23     self.model_B = model_B
24     self.token_pointer = 0
25     self.first = True
```

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 Buffer

func __init__

func estimate_norm_scaling_f...

func refresh

func next

```

26         self.normalize = True
27         self.all_tokens = all_tokens
28
29         estimated_norm_scaling_factor_A = self.estimate_norm_scaling_factor(cfg["mc
30         estimated_norm_scaling_factor_B = self.estimate_norm_scaling_factor(cfg["mc
31
32         self.normalisation_factor = torch.tensor(
33         [
34             estimated_norm_scaling_factor_A,
35             estimated_norm_scaling_factor_B,
36         ],
37         device="cuda:0",
38         dtype=torch.float32,
39         )
40         self.refresh()
41
42     @torch.no_grad()
43     def estimate_norm_scaling_factor(self, batch_size, model, n_batches_for_norm_es
44         # stolen from SAELens https://github.com/jbloomAus/SAELens/blob/6d6eaeef343f
45         norms_per_batch = []
46         for i in tqdm.tqdm(
47             range(n_batches_for_norm_estimate), desc="Estimating norm scaling factc
48         ):
49             tokens = self.all_tokens[i * batch_size : (i + 1) * batch_size]
50             _, cache = model.run_with_cache(
51                 tokens,
52                 names_filter=self.cfg["hook_point"],
53                 return_type=None,
54             )
55             acts = cache[self.cfg["hook_point"]]
56             # TODO: maybe drop BOS here
57             norms_per_batch.append(acts.norm(dim=-1).mean().item())
58         mean_norm = np.mean(norms_per_batch)
59         scaling_factor = np.sqrt(model.cfg.d_model) / mean_norm
60
61         return scaling_factor
62
63     @torch.no_grad()
64     def refresh(self):

```

```

65     self.pointer = 0
66     print("Refreshing the buffer!")
67     with torch.autocast("cuda", torch.bfloat16):
68         if self.first:
69             num_batches = self.buffer_batches
70         else:
71             num_batches = self.buffer_batches // 2
72         self.first = False
73         for _ in tqdm.trange(0, num_batches, self.cfg["model_batch_size"]):
74             tokens = self.all_tokens[
75                 self.token_pointer : min(
76                     self.token_pointer + self.cfg["model_batch_size"], num_batches
77                 )
78             ]
79             _, cache_A = self.model_A.run_with_cache(
80                 tokens, names_filter=self.cfg["hook_point"]
81             )
82             cache_A: ActivationCache
83
84             _, cache_B = self.model_B.run_with_cache(
85                 tokens, names_filter=self.cfg["hook_point"]
86             )
87             cache_B: ActivationCache
88
89             acts = torch.stack([cache_A[self.cfg["hook_point"]], cache_B[self.cfg["hook_point"]]])
90             acts = acts[:, :, 1:, :] # Drop BOS
91             assert acts.shape == (2, tokens.shape[0], tokens.shape[1]-1, self.n_layers, self.d_model)
92             acts = einops.rearrange(
93                 acts,
94                 "n_layers batch seq_len d_model -> (batch seq_len) n_layers d_model"
95             )
96
97             self.buffer[self.pointer : self.pointer + acts.shape[0]] = acts
98             self.pointer += acts.shape[0]
99             self.token_pointer += self.cfg["model_batch_size"]
100
101     self.pointer = 0
102     self.buffer = self.buffer[
103         torch.randperm(self.buffer.shape[0]).to(self.cfg["device"])

```

```
104         ]
105
106     @torch.no_grad()
107     def next(self):
108         out = self.buffer[self.pointer : self.pointer + self.cfg["batch_size"]].float()
109         # out: [batch_size, n_layers, d_model]
110         self.pointer += self.cfg["batch_size"]
111         if self.pointer > self.buffer.shape[0] // 2 - self.cfg["batch_size"]:
112             self.refresh()
113         if self.normalize:
114             out = out * self.normalisation_factor[None, :, None]
115         return out
```