



main

crosscoder-model-diff-replication / crosscoder.py



Go to file

t



ckkissane update default save path

6de8761 · 4 months ago



History

Code

Blame

211 lines (179 loc) · 7.24 KB



Raw



```
1
2 from utils import *
3
4 from torch import nn
5 import pprint
6 import torch.nn.functional as F
7 from typing import Optional, Union
8 from huggingface_hub import hf_hub_download
9
10 from typing import NamedTuple
11
12 DTYPES = {"fp32": torch.float32, "fp16": torch.float16, "bf16": torch.bfloat16}
13 SAVE_DIR = Path("/workspace/crosscoder-model-diff-replication/checkpoints")
14
15 class LossOutput(NamedTuple):
16     # loss: torch.Tensor
17     l2_loss: torch.Tensor
18     l1_loss: torch.Tensor
19     l0_loss: torch.Tensor
20     explained_variance: torch.Tensor
21     explained_variance_A: torch.Tensor
22     explained_variance_B: torch.Tensor
23
24 class CrossCoder(nn.Module):
25     def __init__(self, cfg):
```

Symbols



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



Filter symbols

r

const DTYPES

const SAVE_DIR

class LossOutput

const l2_loss

const l1_loss

const l0_loss

const explained_variance

const explained_variance_A

const explained_variance_B

class CrossCoder

func __init__

func encode

func decode

```

26     super().__init__()
27     self.cfg = cfg
28     d_hidden = self.cfg["dict_size"]
29     d_in = self.cfg["d_in"]
30     self.dtype = DTYPES[self.cfg["enc_dtype"]]
31     torch.manual_seed(self.cfg["seed"])
32     # hardcoding n_models to 2
33     self.W_enc = nn.Parameter(
34         torch.empty(2, d_in, d_hidden, dtype=self.dtype)
35     )
36     self.W_dec = nn.Parameter(
37         torch.nn.init.normal_(
38             torch.empty(
39                 d_hidden, 2, d_in, dtype=self.dtype
40             )
41         )
42     )
43     self.W_dec = nn.Parameter(
44         torch.nn.init.normal_(
45             torch.empty(
46                 d_hidden, 2, d_in, dtype=self.dtype
47             )
48         )
49     )
50     # Make norm of W_dec 0.1 for each column, separate per layer
51     self.W_dec.data = (
52         self.W_dec.data / self.W_dec.data.norm(dim=-1, keepdim=True) * self.cfg
53     )
54     # Initialise W_enc to be the transpose of W_dec
55     self.W_enc.data = einops.rearrange(
56         self.W_dec.data.clone(),
57         "d_hidden n_models d_model -> n_models d_model d_hidden",
58     )
59     self.b_enc = nn.Parameter(torch.zeros(d_hidden, dtype=self.dtype))
60     self.b_dec = nn.Parameter(
61         torch.zeros((2, d_in), dtype=self.dtype)
62     )
63     self.d_hidden = d_hidden
64

```

func forward

func get_losses

func create_save_dir

func save

```

65         self.to(self.cfg["device"])
66         self.save_dir = None
67         self.save_version = 0
68
69     def encode(self, x, apply_relu=True):
70         # x: [batch, n_models, d_model]
71         x_enc = einops.einsum(
72             x,
73             self.W_enc,
74             "batch n_models d_model, n_models d_model d_hidden -> batch d_hidden",
75         )
76         if apply_relu:
77             acts = F.relu(x_enc + self.b_enc)
78         else:
79             acts = x_enc + self.b_enc
80         return acts
81
82     def decode(self, acts):
83         # acts: [batch, d_hidden]
84         acts_dec = einops.einsum(
85             acts,
86             self.W_dec,
87             "batch d_hidden, d_hidden n_models d_model -> batch n_models d_model",
88         )
89         return acts_dec + self.b_dec
90
91     def forward(self, x):
92         # x: [batch, n_models, d_model]
93         acts = self.encode(x)
94         return self.decode(acts)
95
96     def get_losses(self, x):
97         # x: [batch, n_models, d_model]
98         x = x.to(self.dtype)
99         acts = self.encode(x)
100         # acts: [batch, d_hidden]
101         x_reconstruct = self.decode(acts)
102         diff = x_reconstruct.float() - x.float()
103         squared_diff = diff.pow(2)

```

```

104         l2_per_batch = einops.reduce(squared_diff, 'batch n_models d_model -> batch
105         l2_loss = l2_per_batch.mean()
106
107         total_variance = einops.reduce((x - x.mean(0)).pow(2), 'batch n_models d_mc
108         explained_variance = 1 - l2_per_batch / total_variance
109
110         per_token_l2_loss_A = (x_reconstruct[:, 0, :] - x[:, 0, :]).pow(2).sum(dim=
111         total_variance_A = (x[:, 0, :] - x[:, 0, :].mean(0)).pow(2).sum(-1).squeeze
112         explained_variance_A = 1 - per_token_l2_loss_A / total_variance_A
113
114         per_token_l2_loss_B = (x_reconstruct[:, 1, :] - x[:, 1, :]).pow(2).sum(dim=
115         total_variance_B = (x[:, 1, :] - x[:, 1, :].mean(0)).pow(2).sum(-1).squeeze
116         explained_variance_B = 1 - per_token_l2_loss_B / total_variance_B
117
118         decoder_norms = self.W_dec.norm(dim=-1)
119         # decoder_norms: [d_hidden, n_models]
120         total_decoder_norm = einops.reduce(decoder_norms, 'd_hidden n_models -> d_f
121         l1_loss = (acts * total_decoder_norm[None, :]).sum(-1).mean(0)
122
123         l0_loss = (acts>0).float().sum(-1).mean()
124
125         return LossOutput(l2_loss=l2_loss, l1_loss=l1_loss, l0_loss=l0_loss, explai
126
127     def create_save_dir(self):
128         base_dir = Path("/workspace/crosscoder-model-diff-replication/checkpoints")
129         version_list = [
130             int(file.name.split("_")[1])
131             for file in list(SAVE_DIR.iterdir())
132             if "version" in str(file)
133         ]
134         if len(version_list):
135             version = 1 + max(version_list)
136         else:
137             version = 0
138         self.save_dir = base_dir / f"version_{version}"
139         self.save_dir.mkdir(parents=True)
140
141     def save(self):
142         if self.save_dir is None:

```

```

143         self.create_save_dir()
144         weight_path = self.save_dir / f"{self.save_version}.pt"
145         cfg_path = self.save_dir / f"{self.save_version}_cfg.json"
146
147         torch.save(self.state_dict(), weight_path)
148         with open(cfg_path, "w") as f:
149             json.dump(self.cfg, f)
150
151         print(f"Saved as version {self.save_version} in {self.save_dir}")
152         self.save_version += 1
153
154     @classmethod
155     ✓ def load_from_hf(
156         cls,
157         repo_id: str = "ckkissane/crosscoder-gemma-2-2b-model-diff",
158         path: str = "blocks.14.hook_resid_pre",
159         device: Optional[Union[str, torch.device]] = None
160     ) -> "CrossCoder":
161         """
162         Load CrossCoder weights and config from HuggingFace.
163
164         Args:
165             repo_id: HuggingFace repository ID
166             path: Path within the repo to the weights/config
167             model: The transformer model instance needed for initialization
168             device: Device to load the model to (defaults to cfg device if not spec
169
170         Returns:
171             Initialized CrossCoder instance
172         """
173
174         # Download config and weights
175         config_path = hf_hub_download(
176             repo_id=repo_id,
177             filename=f"{path}/cfg.json"
178         )
179         weights_path = hf_hub_download(
180             repo_id=repo_id,
181             filename=f"{path}/cc_weights.pt"

```

```

182     )
183
184     # Load config
185     with open(config_path, 'r') as f:
186         cfg = json.load(f)
187
188     # Override device if specified
189     if device is not None:
190         cfg["device"] = str(device)
191
192     # Initialize CrossCoder with config
193     instance = cls(cfg)
194
195     # Load weights
196     state_dict = torch.load(weights_path, map_location=cfg["device"])
197     instance.load_state_dict(state_dict)
198
199     return instance
200
201     @classmethod
202     ✓ def load(cls, version_dir, checkpoint_version):
203         save_dir = Path("/workspace/crosscoder-model-diff-replication/checkpoints")
204         cfg_path = save_dir / f"{str(checkpoint_version)}_cfg.json"
205         weight_path = save_dir / f"{str(checkpoint_version)}.pt"
206
207         cfg = json.load(open(cfg_path, "r"))
208         pprint.pprint(cfg)
209         self = cls(cfg=cfg)
210         self.load_state_dict(torch.load(weight_path))
211         return self

```