

Code

Blame

59 lines (49 loc) · 1.57 KB

Raw

```
1  # %%
2  from utils import *
3  from crosscoder import CrossCoder
4  torch.set_grad_enabled(False);
5  # %%
6  cross_coder = CrossCoder.load_from_hf()
7
8  # %%
9  norms = cross_coder.W_dec.norm(dim=-1)
10 norms.shape
11 # %%
12 relative_norms = norms[:, 1] / norms.sum(dim=-1)
13 relative_norms.shape
14 # %%
15
16 ✓ fig = px.histogram(
17     relative_norms.detach().cpu().numpy(),
18     title="Gemma 2 2B Base vs IT Model Diff",
19     labels={"value": "Relative decoder norm strength"},
20     nbins=200,
21 )
22
23 fig.update_layout(showlegend=False)
24 fig.update_yaxes(title_text="Number of Latents")
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

r

const

 cross_coder

const

 norms

const

 relative_norms

const

 fig

const

 shared_latent_mask

const

 cosine_sims

const

 fig

```

26 # Update x-axis ticks
27 fig.update_xaxes(
28     tickvals=[0, 0.25, 0.5, 0.75, 1.0],
29     ticktext=['0', '0.25', '0.5', '0.75', '1.0']
30 )
31
32 fig.show()
33
34 # %%
35 shared_latent_mask = (relative_norms < 0.7) & (relative_norms > 0.3)
36 shared_latent_mask.shape
37 # %%
38 # Cosine similarity of recoder vectors between models
39
40 cosine_sims = (cross_coder.W_dec[:, 0, :] * cross_coder.W_dec[:, 1, :]).sum(dim=-1)
41 cosine_sims.shape
42 # %%
43 import plotly.express as px
44 import torch
45
46 ✓ fig = px.histogram(
47     cosine_sims[shared_latent_mask].to(torch.float32).detach().cpu().numpy(),
48     #title="Cosine similarity of decoder vectors between models",
49     log_y=True, # Sets the y-axis to log scale
50     range_x=[-1, 1], # Sets the x-axis range from -1 to 1
51     nbins=100, # Adjust this value to change the number of bins
52     labels={"value": "Cosine similarity of decoder vectors between models"}
53 )
54
55 fig.update_layout(showlegend=False)
56 fig.update_yaxes(title_text="Number of Latents (log scale)")
57
58 fig.show()
59 # %%

```