

 ckkissane first commit

5bde1d5 · 4 months ago

 History

Code

Blame

196 lines (167 loc) · 6.75 KB



```
1  # %%
2  import os
3  from IPython import get_ipython
4
5  ipython = get_ipython()
6  # Code to automatically update the HookedTransformer code as its edited without res
7  if ipython is not None:
8      ipython.magic("load_ext autoreload")
9      ipython.magic("autoreload 2")
10
11  import plotly.io as pio
12  pio.renderers.default = "jupyterlab"
13
14  # Import stuff
15  import einops
16  import json
17  import argparse
18
19  from datasets import load_dataset
20  from pathlib import Path
21  import plotly.express as px
22  from torch.distributions.categorical import Categorical
23  from tqdm import tqdm
24  import torch
25  import numpy as np
```

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 ipython

const update_layout_set

func imshow

func line

func scatter

func lines

func bar

func create_html

func arg_parse_update_cfg

func load_pile_lmsys_mixed_t...

```

26 from transformer_lens import HookedTransformer
27 from jaxtyping import Float
28 from transformer_lens.hook_points import HookPoint
29
30 from functools import partial
31
32 from IPython.display import HTML
33
34 from transformer_lens.utils import to_numpy
35 import pandas as pd
36
37 from html import escape
38 import colorsys
39
40
41 import wandb
42
43 import plotly.graph_objects as go
44
45 ✓ update_layout_set = {
46     "xaxis_range", "yaxis_range", "hovermode", "xaxis_title", "yaxis_title", "color
47     "title_x", "bargap", "bargroupgap", "xaxis_tickformat", "yaxis_tickformat", "t
48     "xaxis_gridwidth", "xaxis_gridcolor", "yaxis_showgrid", "yaxis_gridwidth"
49 }
50
51 ✓ def imshow(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
52     if isinstance(tensor, list):
53         tensor = torch.stack(tensor)
54     kwargs_post = {k: v for k, v in kwargs.items() if k in update_layout_set}
55     kwargs_pre = {k: v for k, v in kwargs.items() if k not in update_layout_set}
56     if "facet_labels" in kwargs_pre:
57         facet_labels = kwargs_pre.pop("facet_labels")
58     else:
59         facet_labels = None
60     if "color_continuous_scale" not in kwargs_pre:
61         kwargs_pre["color_continuous_scale"] = "RdBu"
62     fig = px.imshow(to_numpy(tensor), color_continuous_midpoint=0.0, labels={"x": xax
63     if facet_labels:
64         for i, label in enumerate(facet_labels):

```

```

65         fig.layout.annotations[i]['text'] = label
66
67     fig.show(renderer)
68
69     def line(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
70         px.line(y=to_numpy(tensor), labels={"x":xaxis, "y":yaxis}, **kwargs).show(rende
71
72     def scatter(x, y, xaxis="", yaxis="", caxis="", renderer=None, return_fig=False, **
73         x = to_numpy(x)
74         y = to_numpy(y)
75         fig = px.scatter(y=y, x=x, labels={"x":xaxis, "y":yaxis, "color":caxis}, **kwar
76         if return_fig:
77             return fig
78         fig.show(renderer)
79
80     def lines(lines_list, x=None, mode='lines', labels=None, xaxis='', yaxis='', title
81         # Helper function to plot multiple lines
82         if type(lines_list)==torch.Tensor:
83             lines_list = [lines_list[i] for i in range(lines_list.shape[0])]
84         if x is None:
85             x=np.arange(len(lines_list[0]))
86         fig = go.Figure(layout={'title':title})
87         fig.update_xaxes(title=xaxis)
88         fig.update_yaxes(title=yaxis)
89         for c, line in enumerate(lines_list):
90             if type(line)==torch.Tensor:
91                 line = to_numpy(line)
92             if labels is not None:
93                 label = labels[c]
94             else:
95                 label = c
96             fig.add_trace(go.Scatter(x=x, y=line, mode=mode, name=label, hovertext=hove
97         if log_y:
98             fig.update_layout(yaxis_type="log")
99         fig.show()
100
101     def bar(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
102         px.bar(
103             y=to_numpy(tensor),

```

```

104         labels={"x": xaxis, "y": yaxis},
105         template="simple_white",
106         **kwargs).show(renderer)
107
108  ✓ def create_html(strings, values, saturation=0.5, allow_different_length=False):
109     # escape strings to deal with tabs, newlines, etc.
110     escaped_strings = [escape(s, quote=True) for s in strings]
111     processed_strings = [
112         s.replace("\n", "<br/>").replace("\t", "&emsp;").replace(" ", "&nbsp;")
113         for s in escaped_strings
114     ]
115
116     if isinstance(values, torch.Tensor) and len(values.shape)>1:
117         values = values.flatten().tolist()
118
119     if not allow_different_length:
120         assert len(processed_strings) == len(values)
121
122     # scale values
123     max_value = max(max(values), -min(values))+1e-3
124     scaled_values = [v / max_value * saturation for v in values]
125
126     # create html
127     html = ""
128     for i, s in enumerate(processed_strings):
129         if i<len(scaled_values):
130             v = scaled_values[i]
131         else:
132             v = 0
133         if v < 0:
134             hue = 0 # hue for red in HSV
135         else:
136             hue = 0.66 # hue for blue in HSV
137         rgb_color = colorsys.hsv_to_rgb(
138             hue, v, 1
139         ) # hsv color with hue 0.66 (blue), saturation as v, value 1
140         hex_color = "#%02x%02x%02x" % (
141             int(rgb_color[0] * 255),
142             int(rgb_color[1] * 255),

```

```

143         int(rgb_color[2] * 255),
144     )
145     html += f'<span style="background-color: {hex_color}; border: 1px solid lig
146
147     display(HTML(html))
148
149     # crosscoder stuff
150
151     ✓ def arg_parse_update_cfg(default_cfg):
152         """
153         Helper function to take in a dictionary of arguments, convert these to command
154
155         If in Ipython, just returns with no changes
156         """
157         if get_ipython() is not None:
158             # Is in IPython
159             print("In IPython - skipped argparse")
160             return default_cfg
161         cfg = dict(default_cfg)
162         parser = argparse.ArgumentParser()
163         for key, value in default_cfg.items():
164             if type(value) == bool:
165                 # argparse for Booleans is broken rip. Now you put in a flag to change
166                 if value:
167                     parser.add_argument(f"--{key}", action="store_false")
168                 else:
169                     parser.add_argument(f"--{key}", action="store_true")
170
171             else:
172                 parser.add_argument(f"--{key}", type=type(value), default=value)
173         args = parser.parse_args()
174         parsed_args = vars(args)
175         cfg.update(parsed_args)
176         print("Updated config")
177         print(json.dumps(cfg, indent=2))
178         return cfg
179
180     ✓ def load_pile_lmsys_mixed_tokens():
181         try:

```

```
182     print("Loading data from disk")
183     all_tokens = torch.load("/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-
184 except:
185     print("Data is not cached. Loading data from HF")
186     data = load_dataset(
187         "ckkissane/pile-lmsys-mix-1m-tokenized-gemma-2",
188         split="train",
189         cache_dir="/workspace/cache/"
190     )
191     data.save_to_disk("/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-2.hf")
192     data.set_format(type="torch", columns=["input_ids"])
193     all_tokens = data["input_ids"]
194     torch.save(all_tokens, "/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-2
195     print(f"Saved tokens to disk")
196     return all_tokens
```