

First diff is for MBRL-Lib

Second diff is for Gym-Duckietown

## Comparing changes

This is a direct comparison between two commits made in this repository or its related repositories. View the default comparison for this range [here](#).

base: c213d0b ▾ ↔ compare: 1fc1e15 ▾

Showing 20 changed files with 2,546 additions and 14 deletions.

[Split](#) [Unified](#)

1 .gitignore

28	28	MANIFEST
29	29	wandb/
30	30	scratch.py
	31	+ scratch.txt

42 README.md

28	28	Then to activate it use:
29	29	
30	30	conda activate RLDucky
	31	+
	32	+ You will use this python virtual environment whenever you are running code from this repo. You should use a separate virtual environment when running code from the duckietown repo, such as in the commands below.
31	33	### Gym-Duckietown
32	34	<!-- You will need to do the [duckietown laptop] ( <a href="https://docs.duckietown.org/daffy/opmanual_duckiebot/out/laptop_setup.html">https://docs.duckietown.org/daffy/opmanual_duckiebot/out/laptop_setup.html</a> ) setup to use the gym-duckietown -->
	35	+
	36	+ Make a new virtual environment for the Gym-Duckietown repo:
	37	+
	38	conda create --name GymDucky python=3.8
	39	+
	40	Then to activate it use:
	41	+
	42	conda activate GymDucky
	43	+
	44	+
33	45	The first repo to clone is [ <a href="#">Gym-Duckietown</a> ]( <a href="https://github.com/duckietown/gym-duckietown.git">https://github.com/duckietown/gym-duckietown.git</a> ), and make sure you checkout and use the master branch. Additionally you can install the required python packages for that repo via the command `pip install -e .` where `.` specifies the current directory.
34	46	
	47	+     conda activate GymDucky
35	48	git clone <a href="https://github.com/duckietown/gym-duckietown.git">https://github.com/duckietown/gym-duckietown.git</a>
36	49	cd gym-duckietown
37	50	git checkout master
42	55	While trying to use Gym-Duckietown we ran into an issue involving a malfunctioning / deprecated `geometry` module. If you run into the same problem, you can just comment out that import. So just navigate to the `gym-duckietown/gym_duckietown/simulator.py` file and comment out the `import geometry` line.
43	56	
44	57	### Importing Gym-Duckietown into MBRL-Lib
45	58	-
46	59	+ For the following commands, switch back to your RLDucky environment:
	60	+
	61	conda activate RLDucky
	62	+
47	63	To use the Duckietown environment seamlessly with MBRL-Lib, we will have to add the `gym-duckietown` repo as a python module to our python installation. There are two ways of doing this.
48	64	
49	65	#### Option 1: Using Path (.pth) Files

```

100 116     Clone this repository and install the required python packages:
101 117
102 118         git clone https://github.com/alik-git/duckietown-mbrl-lib
103 119 - cd mbrl-lib
104 120 + cd duckietown-mbrl-lib
105 121     conda activate RLucky
106 122     pip install -e ".[ducky]"
107 123
108 124     cd /your/path/to/mbrl-lib
109 125     python -m pytest tests/mujoco
110 126
111 127
112 128
113 129
114 130 - ##### Side Note:
115 131 - While trying to run MuJoCo we twice ran into an error relating to something involving `undefined symbol: __glewBindBuffer`, and the only fix we found (from a Reddit [thread] (https://www.reddit.com/r/reinforcementlearning/comments/qay11a/how\_to\_use\_mujoco\_from\_python3/)) was to
116 132     install the following packages:
117 133 + ##### Mujoco Issues:
118 134 +
119 135 + A few issues with mujoco - and I've encountered this during other projects as well - are about mujoco not
120 136     neatly attaching to a rendering engine in your software environment. On Ubuntu it seems to try and use the
121 137     [OpenGL] (https://en.wikipedia.org/wiki/OpenGL) rendering engine, and to find it it uses [GLEW]
122 138     (https://github.com/nigels-com/glew). But this process doesn't always end up being seamless. Below are a few
123 139     issues that we ran into:
124 140 + ##### Rendering Issue 1:
125 141 +
126 142 + While trying to run MuJoCo we ran into an error twice relating to something involving `undefined symbol:
127 143     __glewBindBuffer`, and the only fix we found (from a Reddit [thread]
128 144     (https://www.reddit.com/r/reinforcementlearning/comments/qay11a/how\_to\_use\_mujoco\_from\_python3/)) was to
129 145     install the following packages:
130 146
131 147     sudo apt install curl git libgl1-mesa-dev libgl1-mesa-glx libglew-dev \
132 148     libosmesa6-dev software-properties-common net-tools unzip vim \
133 149     virtualenv wget xpra xserver-xorg-dev libglfw3-dev patchelf
134 150
135 151 + ##### Rendering Issue 2:
136 152 +
137 153 + This issue is again about the [dm_control] (https://github.com/deepmind/dm\_control) installation of mujoco, you
138 154     basically get an [error] (https://github.com/deepmind/dm\_control/issues/283) saying `mujoco.FatalError:
139 155     gladLoadGL error`, and as described in the [aforementioned link]
140 156     (https://github.com/deepmind/dm\_control/issues/283#issuecomment-1095490151), the solution is to find your
141 157     dm_control package in your python installation, something like:
142 158 +
143 159 + ``
144 160 + /home/username/anaconda3/envs/RLucky/lib/python3.8/site-packages/dm\_control/\_render/glfw\_renderer.py
145 161 +
146 162 + and modify it to add some extra lines as described [here]
147 163     (https://github.com/deepmind/dm\_control/issues/283#issuecomment-1095490151).
148 164 +
149 165 #### Logging and Visualization (W&B)
150 166
151 167 We use [Weights & Biases] (https://wandb.ai/site) for logging and visualizing our run metrics. If you're
152 168 unfamiliar with Weights & Biases, it is a powerful and convenient library to organize and track ML experiments.
153 169 You can take look at their [quick-start guide] (https://docs.wandb.ai/quickstart) and [documentation]
154 170     (https://docs.wandb.ai/), and you'll have to create an account to be able to view and use the dashboard, you
155 171     can do so [here] (https://wandb.ai/site).
156 172
157 173 To run an experiment you can use commands in the following format:
158 174
159 175     python -m mbrl.examples.main algorithm=planet dynamics_model=planet overrides=planet_duckietown
160 176
161 177 +
162 178 + Here is another example of a command with shorter episodes:
163 179 +
164 180     python -m mbrl.examples.main algorithm=planet dynamics_model=planet overrides=planet_cheetah_run
165 181     algorithm.test_frequency=2 overrides.sequence_length=10 overrides.batch_size=10
166 182
167 183
168 184
169 185
170 186
171 187 You will see the output of your run in the terminal as well as in a results file created by Hydra located by
172 188 default at `.

```

(`./exp`) by passing  
`root\_dir=path-to-your-dir`, and the experiment sub-folder (`default`) by

451 base\_sweep\_config.yaml □

Load diff

Large diffs are not rendered by default.

300 mbrl/algorithms/dreamer.py □

```
... ... @@ -0,0 +1,300 @@
1 + ##### FOR GLEN #####
2 + #####
3 + #####
4 + # This file is basically a copy of MBRL-Lib's PlaNet, modified to accommodate
5 + # the dreamer method.
6 + #####
7 +
8 + # Copyright (c) Facebook, Inc. and its affiliates. All Rights Reserved.
9 + #
10 + # This source code is licensed under the MIT license found in the
11 + # LICENSE file in the root directory of this source tree.
12 + import os
13 + import pathlib
14 + from typing import List, Optional, Union, cast
15 +
16 + import gym
17 + import hydra
18 + import numpy as np
19 + import omegaconf
20 + import torch
21 +
22 + import mbrl.constants
23 + #import mbrl.third_party.pytorch_sac as pytorch_sac
24 +
25 + from mbrl.env.termination_fns import no_termination
26 + from mbrl.models import ModelEnv, ModelTrainer
27 + from mbrl.planning import RandomAgent, create_trajectory_optim_agent_for_model
28 + from mbrl.util import Logger
29 + from mbrl.util.common import (
30 +     create_replay_buffer,
31 +     get_sequence_buffer_iterator,
32 +     rollout_agent_trajectories,
33 + )
34 + #from mbrl.planning.sac_wrapper import SACAgent
35 + from mbrl.planning.dreamer_wrapper import DreamerAgent
36 +
37 + import wandb
38 + from gym.wrappers import Monitor
39 +
40 +
41 + # Original modified from PlaNet
42 +
43 + METRICS_LOG_FORMAT = [
44 +     ("observations_loss", "OL", "float"),
45 +     ("reward_loss", "RL", "float"),
46 +     ("gradient_norm", "GN", "float"),
47 +     ("kl_loss", "KL", "float"),
48 + ]
49 +
50 +
51 + def train(
52 +     env: gym.Env,
53 +     cfg: omegaconf.DictConfig,
```

```

54 +         silent: bool = False,
55 +         work_dir: Union[Optional[str], pathlib.Path] = None,
56 +     ) -> np.float32:
57 +         # Experiment initialization
58 +         debug_mode = cfg.get("debug_mode", False)
59 +
60 +         if work_dir is None:
61 +             work_dir = os.getcwd()
62 +             work_dir = pathlib.Path(work_dir)
63 +             print(f"Results will be saved at {work_dir}.")
64 +             wandb.config.update({"work_dir": str(work_dir)})
65 +
66 +         if silent:
67 +             logger = None
68 +         else:
69 +             logger = Logger(work_dir)
70 +             logger.register_group("metrics", METRICS_LOG_FORMAT, color="yellow")
71 +             logger.register_group(
72 +                 mbirl.constants.RESULTS_LOG_NAME,
73 +                 [
74 +                     ("env_step", "S", "int"),
75 +                     ("train_episode_reward", "RT", "float"),
76 +                     ("episode_reward", "ET", "float"),
77 +                 ],
78 +                 color="green",
79 +             )
80 +
81 +         rng = torch.Generator(device=cfg.device)
82 +         rng.manual_seed(cfg.seed)
83 +         np_rng = np.random.default_rng(seed=cfg.seed)
84 +
85 +         # Create replay buffer and collect initial data
86 +         replay_buffer = create_replay_buffer(
87 +             cfg,
88 +             env.observation_space.shape,
89 +             env.action_space.shape,
90 +             collect_trajectories=True,
91 +             rng=np_rng,
92 +         )
93 +         rollout_agent_trajectories(
94 +             env,
95 +             cfg.algorithm.num_initial_trajectories,
96 +             RandomAgent(env),
97 +             agent_kwargs={},
98 +             replay_buffer=replay_buffer,
99 +             collect_full_trajectories=True,
100 +             trial_length=cfg.overrides.trial_length,
101 +             agent_uses_low_dim_obs=False,
102 +         )
103 +
104 +         # Create PlaNet model
105 +         cfg.dynamics_model.action_size = env.action_space.shape[0]
106 +
107 +         # Use hydra to create a dreamer model (really uses PlaNet model)
108 +         dreamer = hydra.utils.instantiate(cfg.dynamics_model)
109 +         # Give it the real gym env to model
110 +         dreamer.setGymEnv(env, work_dir)
111 +
112 +         # adam optim that takes into account all 3 network losses
113 +         # actor, critic, model
114 +         dreamer_optim = dreamer.configure_optimizers()
115 +         assert isinstance(dreamer, mbirl.models.DreamerModel)
116 +         model_env = ModelEnv(env, dreamer, no_termination, generator=rng)
117 +         trainer = ModelTrainer(dreamer, logger=logger, optim_lr=1e-3, optim_eps=1e-4)
118 +
119 +         #####
120 +         ##### FOR GLEN #####
121

```

```

122 + #####
123 + #
124 +
125 + # This library has a function called
126 + # create_trajectory_optim_agent_for_model() that is used to create an agent
127 + # (an agent that essentially uses CEM planning) given a world model.
128 +
129 + #
130 + # Example: agent =
131 + # create_trajectory_optim_agent_for_model(model_env, cfg.algorithm.agent)
132 +
133 + # But in our case, dreamer is the world model i.e. model_env, but Dreamer
134 + # also includes an "agent" in the form of an actor and critic network, and
135 + # Dreamer therefore doesn't need CEM planning.
136 +
137 + # So instead of doing it the "MBRL-Lib" way and creating an agent, we just
138 + # use the world model i.e. Dreamer to sample actions from. This doesn't sit
139 + # well with the structure of this library so in future work we will try to
140 + # separate the model part of Dreamer (which is just PlaNet) and the
141 + # agent/planning part of Dreamer, so that it fits better with this library.
142 +
143 +
144 + # Create Dreamer Agent (Action and Value model), are these needed for this to operate properly?
145 + # This agent rolls out trajectories using ModelEnv, which uses planet.sample()
146 + # to simulate the trajectories from the prior transition model
147 + # The starting point for trajectories is conditioned on the latest observation,
148 + # for which we use planet.update_posterior() after each environment step
149 + # the CEM way
150 + # agent = create_trajectory_optim_agent_for_model(model_env, cfg.algorithm.agent)
151 + # the SAC/Dreamer way
152 + # agent = different_or_same_function(model_env, cfg.algorithm.agent)
153 +
154 +
155 + # Callback and containers to accumulate training statistics and average over batch
156 + rec_losses: List[float] = []
157 + reward_losses: List[float] = []
158 + kl_losses: List[float] = []
159 + grad_norms: List[float] = []
160 +
161 + def get_metrics_and_clear_metric_containers():
162 +     metrics_ = {
163 +         "observations_loss": np.mean(rec_losses).item(),
164 +         "reward_loss": np.mean(reward_losses).item(),
165 +         "gradient_norm": np.mean(grad_norms).item(),
166 +         "kl_loss": np.mean(kl_losses).item(),
167 +     }
168 +
169 +     for c in [rec_losses, reward_losses, kl_losses, grad_norms]:
170 +         c.clear()
171 +
172 +     return metrics_
173 +
174 + def batch_callback(_epoch, _loss, meta, _mode):
175 +     if meta:
176 +         rec_losses.append(meta["observations_loss"])
177 +         reward_losses.append(meta["reward_loss"])
178 +         kl_losses.append(meta["kl_loss"])
179 +         if "grad_norm" in meta:
180 +             grad_norms.append(meta["grad_norm"])
181 +
182 + def is_test_episode(episode_):
183 +     return episode_ % cfg.algorithm.test_frequency == 0
184 +
185 + # PlaNet loop
186 + step = replay_buffer.num_stored
187 + total_rewards = 0.0
188 + for episode in range(cfg.algorithm.num_episodes):
189

```

```

189 +         dreamer.set_curr_episode(episode)
190 +
191 +     # Train the model for one epoch of `num_grad_updates`
192 +     dataset, _ = get_sequence_buffer_iterator(
193 +         replay_buffer,
194 +         cfg.overrides.batch_size,
195 +         0, # no validation data
196 +         cfg.overrides.sequence_length,
197 +         max_batches_per_loop_train=cfg.overrides.num_grad_updates,
198 +         use_simple_sampler=True,
199 +     )
200 +     trainer.train(
201 +         dataset, num_epochs=1, batch_callback=batch_callback, evaluate=False
202 +     )
203 +     dreamer.save(work_dir / "dreamer.pth")
204 +     replay_buffer.save(work_dir)
205 +     metrics = get_metrics_and_clear_metric_containers()
206 +     logger.log_data("metrics", metrics)
207 +     wandb_metrics = metrics
208 +     wandb_metrics["global_episode"] = episode
209 +     wandb.log(wandb_metrics)
210 +
211 +     if is_test_episode(episode):
212 +         print("AHH ITS A TEST EPISODE!!!")
213 +         curr_env = Monitor(env, work_dir, force=True)
214 +         dreamer.set_currently_testing(True)
215 +     else:
216 +         curr_env = env
217 +         dreamer.set_currently_testing(False)
218 +
219 +     # Collect one episode of data
220 +     episode_reward = 0.0
221 +     obs = curr_env.reset()
222 +     # want to do
223 +     #agent.reset()
224 +     dreamer.reset_world_model(device=cfg.device)
225 +     state = None
226 +     action = None
227 +     done = False
228 +     while not done:
229 +
230 +         # hacky check to see if we are using the duckietown environment
231 +         # I just see if the obs shape ends in 3, for the other envs it does not
232 +         if obs.shape[-1] == 3:
233 +             dreamer.in_duckietown = True
234 +
235 +             if dreamer.in_duckietown:
236 +                 obs = np.transpose(obs, (2,0,1))
237 +                 pass
238 +
239 +             # want to do
240 +             # func dreamer.update(...)
241 +             # planet.update(..)
242 +             # actor_net.update(..)
243 +             # value_net.update(..)
244 +             # dreamer.update_alg(obs, action=action, rng=rng)
245 +             ...
246 +             Don't need yet, noise in implementation
247 +             action_noise = (
248 +                 0
249 +                 if is_test_episode(episode)
250 +                 else cfg.overrides.action_noise_std
251 +                 * np_rng.standard_normal(curr_env.action_space.shape[0])
252 +             )
253 +             ...
254 +             # want to do (kinda)
255 +             # action, _ = dreamer.policy(obs)

```

```

256 +         #action sampler for Dreamer wants tensors with a batch dimension
257 +         action, _, state, = dreamer.action_sampler_fn(torch.FloatTensor(obs).unsqueeze(0), state, dreamer.ex
258 +         ''
259 +         Already have noise in the implementation
260 +         action = action + action_noise
261 +         # action = agent.act(obs) + action_noise
262 +         action = np.clip(action, -1.0, 1.0) # to account for the noise
263 +         ''
264 +
265 +         if dreamer.in_duckietown:
266 +             obs = np.transpose(obs, (1,2,0))
267 +             action = action.squeeze(0).cpu().numpy() #env likes numpy arrays for actions unbatched
268 +             next_obs, reward, done, info = curr_env.step(action)
269 +             replay_buffer.add(obs, action, next_obs, reward, done)
270 +             episode_reward += reward
271 +             obs = next_obs
272 +             if debug_mode:
273 +                 print(f"step: {step}, reward: {reward}.")
274 +             step += 1
275 +             total_rewards += episode_reward
276 +             logger.log_data(
277 +                 mbrl.constants.RESULTS_LOG_NAME,
278 +                 {
279 +                     "episode_reward": episode_reward * is_test_episode(episodes),
280 +                     "train_episode_reward": episode_reward * (1 - is_test_episode(episodes)),
281 +                     "env_step": step,
282 +                 },
283 +             )
284 +             wandb.log(
285 +                 {
286 +                     "episode_reward": episode_reward * is_test_episode(episodes),
287 +                     "train_episode_reward": episode_reward * (1 - is_test_episode(episodes)),
288 +                     "env_step": step,
289 +                     "global_episode": episodes
290 +                 }
291 +             )
292 +             avg_ep_reward = total_rewards / (episodes+1)
293 +             wandb.log({'average_episode_reward': avg_ep_reward, "global_episode": episodes})
294 +
295 +             # returns average episode reward (e.g., to use for tuning learning curves)
296 +             avg_ep_reward = total_rewards / cfg.algorithm.num_episodes
297 +             wandb.log({'average_episode_reward': avg_ep_reward, "global_episode": episodes})
298 +             return avg_ep_reward
299 +
300 +

```

▼ ⏪ 40 mbrl/algorithms/planet.py □

...	...	@@ -1,3 +1,8 @@
	1	+ #####
	2	+ ##### FOR GLEN #####
	3	+ #####
	4	+ # We added video recording for the Duckietown environments and better logging.
	5	+ #####
1	6	# Copyright (c) Facebook, Inc. and its affiliates. All Rights Reserved.
2	7	#
3	8	# This source code is licensed under the MIT license found in the
24	29	)
25	30	
26	31	import wandb
	32	+ from gym.wrappers import Monitor
	33	+ from gym.wrappers.monitoring.video_recorder import VideoRecorder
	34	+ from PIL import Image
	35	+ import cv2
27	36	
28	37	METRICS_LOG_FORMAT = [

```

29   38     ("observations_loss", "OL", "float"),
46   55         work_dir = os.getcwd()
47   56     work_dir = pathlib.Path(work_dir)
48   57     print(f"Results will be saved at {work_dir}.")
58 +     wandb.config.update({"work_dir": str(work_dir)})
49   59
50   60     if silent:
51   61         logger = None
90  100     planet = hydra.utils.instantiate(cfg.dynamics_model)
91  101     assert isinstance(planet, mbrl.models.PlaNetModel)
92  102     model_env = ModelEnv(env, planet, no_termination, generator=rng)
93 -     trainer = ModelTrainer(planet, logger=logger, optim_lr=1e-3, optim_eps=1e-4)
103 +     print(f"\n###\nUsing learning rate: {cfg.overrides.model_learning_rate}\n###")
104 +
105 +     trainer = ModelTrainer(planet, logger=logger, optim_lr=cfg.overrides.model_learning_rate, optim_eps=1e-4)
94  106
95  107     # Create CEM agent
96  108     # This agent rolls outs trajectories using ModelEnv, which uses planet.sample()
133 145     step = replay_buffer.num_stored
134 146     total_rewards = 0.0
135 147     for episode in range(cfg.algorithm.num_episodes):
148 +     print(f"##\nNow on episode {episode} of {cfg.algorithm.num_episodes - 1}")
136 149     # Train the model for one epoch of `num_grad_updates`
137 150     dataset, _ = get_sequence_buffer_iterator(
138 151         replay_buffer,
148 161         planet.save(work_dir / "planet.pth")
149 162         replay_buffer.save(work_dir)
150 163         metrics = get_metrics_and_clear_metric_containers()
164 +     wandb_metrics = metrics
165 +     wandb_metrics['global_episode'] = episode
151 166     logger.log_data("metrics", metrics)
152 -     wandb.log(metrics)
167 +
168 +     wandb.log(wandb_metrics)
169 +
170 +     if is_test_episode(episode):
171 +         print(f"Reached Test Episode! Episode: {episode}")
172 +         curr_env = Monitor(env, work_dir / "videos" / f"{episode:05d}", force=True)
173 +         planet.save(work_dir / "models" / f"{episode:05d}" / "planet.pth")
174 +
175 +     else:
176         curr_env = env
153 176
154 177     # Collect one episode of data
155 178     episode_reward = 0.0
156 -     obs = env.reset()
179 +     obs = curr_env.reset()
157 180     agent.reset()
158 181     planet.reset_posterior()
159 182     action = None
164 187         0
165 188             if is_test_episode(episode)
166 189             else cfg.overrides.action_noise_std
167 -             * np_rng.standard_normal(env.action_space.shape[0])
190 +             * np_rng.standard_normal(curr_env.action_space.shape[0])
168 191         )
169 192         action = agent.act(obs) + action_noise
170 193         action = np.clip(action, -1.0, 1.0) # to account for the noise
171 -         next_obs, reward, done, info = env.step(action)
194 +         next_obs, reward, done, info = curr_env.step(action)
172 195         replay_buffer.add(obs, action, next_obs, reward, done)
173 196         episode_reward += reward
174 197         obs = next_obs
189 212             "episode_reward": episode_reward * is_test_episode(episode),
190 213             "train_episode_reward": episode_reward * (1 - is_test_episode(episode)),
191 214             "env_step": step,
215 +             "global_episode": episode
192 216         }
193 217

```

```

        )
218 +     avg_ep_reward = total_rewards / (episode+1)
219 +     wandb.log({'average_episode_reward': avg_ep_reward, "global_episode": episode})
194 220
195 221     # returns average episode reward (e.g., to use for tuning learning curves)
196 222 -     return total_rewards / cfg.algorithm.num_episodes
222 +     avg_ep_reward = total_rewards / cfg.algorithm.num_episodes
223 +     wandb.log({'average_episode_reward': avg_ep_reward, "global_episode": episode})
224 +     return avg_ep_reward

```

▼ 59 mbrl/examples/conf/algorithm/dreamer.yaml □

```

... ... @@ -0,0 +1,59 @@
1 + # @package _group_
2 +
3 + #####
4 + ##### FOR GLEN #####
5 + #####
6 + # This is the config file for the "planning parts" of Dreamer
7 + #####
8 +
9 + name: "dreamer"
10 +
11 + agent:
12 +     # this calls a Dreamer agent that we made
13 +     _target_: mbrl.planning.dreamer_wrapper.DreamerAgent #mbrl.third_party.pytorch_sac.agent.sac.SACAgent
14 +     env: ${overrides.environment}
15 +
16 + num_initial_trajectories: 5
17 + action_noise_std: 0.3
18 + test_frequency: 25
19 + num_episodes: 1000
20 + dataset_size: 1000000
21 +
22 +     # these are from an SAC implementation that we may use
23 +     # from mbrl.planning.sac_wrapper import SACAgent
24 +     # obs_dim: ??? # to be specified later
25 +     # action_dim: ??? # to be specified later
26 +     # action_range: ??? # to be specified later
27 +     # device: ${device}
28 +
29 +     #     critic_cfg: ${algorithm.double_q_critic}
30 +     #     actor_cfg: ${algorithm.diag_gaussian_actor}
31 +     #     discount: 0.99
32 +     #     init_temperature: 0.1
33 +     #     alpha_lr: ${overrides.sac_alpha_lr}
34 +     #     alpha_betas: [0.9, 0.999]
35 +     #     actor_lr: ${overrides.sac_actor_lr}
36 +     #     actor_betas: [0.9, 0.999]
37 +     #     actor_update_frequency: ${overrides.sac_actor_update_frequency}
38 +     #     critic_lr: ${overrides.sac_critic_lr}
39 +     #     critic_betas: [0.9, 0.999]
40 +     #     critic_tau: 0.005
41 +     #     critic_target_update_frequency: ${overrides.sac_critic_target_update_frequency}
42 +     #     batch_size: 256
43 +     #     learnable_temperature: true
44 +     #     target_entropy: ${overrides.sac_target_entropy}
45 +
46 +     # double_q_critic:
47 +     #     _target_: mbrl.third_party.pytorch_sac.agent.critic.DoubleQCritic
48 +     #     obs_dim: ${algorithm.agent.obs_dim}
49 +     #     action_dim: ${algorithm.agent.action_dim}
50 +     #     hidden_dim: 1024
51 +     #     hidden_depth: ${overrides.sac_hidden_depth}
52 +
53 +     # diag_gaussian_actor:
54 +     #     _target_: mbrl.third_party.pytorch_sac.agent.actor.DiagGaussianActor
55 +

```

```

56     #   obs_dim: ${algorithm.agent.obs_dim}
57 + #   action_dim: ${algorithm.agent.action_dim}
58 + #   hidden_depth: ${overrides.sac_hidden_depth}
59 + #   hidden_dim: 1024
59 + #   log_std_bounds: [-5, 2] ⊖

```

⌄ 37 mbrl/examples/conf/dynamics\_model/dreamer.yaml □

```

... ... @@ -0,0 +1,37 @@
1 + # @package _group_
2 + #####
3 + ##### FOR GLEN #####
4 + #####
5 + # This is the config file for the "model parts" of Dreamer. Again this mostly
6 + # borrows settings from the MBRL-Lib PlaNet and the recommended config from
7 + # Chandramouli Rajagopalan's pytorch implementation of Dreamer.
8 + #####
9 + _target_: mbrl.models.DreamerModel
10 + obs_shape: [3, 64, 64]
11 + # obs_encoding_size: 1024
12 + # encoder_config:
13 + #   - [3, 32, 4, 2]
14 + #   - [32, 64, 4, 2]
15 + #   - [64, 128, 4, 2]
16 + #   - [128, 256, 4, 2]
17 + # decoder_config:
18 + #   - [1024, 1, 1]
19 + #   - [1024, 128, 5, 2]
20 + #   - [128, 64, 5, 2]
21 + #   - [64, 32, 6, 2]
22 + #   - [32, 3, 6, 2]
23 + action_size: ???
24 + hidden_size_fcs: 400
25 + depth_size: 32
26 + stoch_size: 30
27 + deter_size: 200
28 + # belief_size: 200
29 + # latent_state_size: 30
30 + device: ${device}
31 + min_std: 0.1
32 + # free_nats: 3.0
33 + # kl_scale: 1.0
34 + # grad_clip_norm: 1000.0
35 + # # To use for duckietown_gym_env
36 + # grad_clip_norm: 10.0
37 + # env: ${overrides.env}

```

⌄ 42 mbrl/examples/conf/overrides/dreamer\_cheetah\_run.yaml □

```

... ... @@ -0,0 +1,42 @@
1 + # @package _group_
2 + #####
3 + ##### FOR GLEN #####
4 + #####
5 + # This is the config file to run Dreamer on Cheetah Run. For now we mostly just
6 + # use the same values as PlaNet.
7 + #####
8 + env: "dmcontrol_cheetah_run" # used to set the hydra dir, ignored otherwise
9 +
10 + env_cfg:
11 + _target_: "mbrl.third_party.dmc2gym.wrappers.DMCWrapper"
12 + domain_name: "cheetah"
13 + task_name: "run"
14 + task_kwargs:
15 +   random: ${seed}
16 +   visualize_reward: false
17 +   from_pixels: true
18 +   height: 64

```

```

19 +   width: 64
20 +   frame_skip: 4
21 +   bit_depth: 5
22 +
23 +   term_fn: "no_termination"
24 +
25 + # General configuration overrides
26 + trial_length: 250
27 + action_noise_std: 0.3
28 +
29 + # Model overrides
30 + num_grad_updates: 100
31 + sequence_length: 50
32 + batch_size: 50
33 + free_nats: 3
34 + kl_scale: 1.0
35 +
36 + # Planner configuration overrides
37 + planning_horizon: 12
38 + cem_num_iters: 10
39 + cem_elite_ratio: 0.1
40 + cem_population_size: 1000
41 + cem_alpha: 0.0
42 + cem_clipped_normal: true

```

✓ 32 mbrl/examples/conf/overrides/dreamer\_duckietown.yaml ↗

...	...	@@ -0,0 +1,32 @@
	1	+ # @package _group_
	2	+ #####
	3	+ ##### FOR GLEN #####
	4	+ #####
	5	+ # This is the config file to run Dreamer on Duckietown. For now we mostly just
	6	+ # use the same values as PlaNet.
	7	+ #####
	8	+ env: "duckietown_gym_env" # used to set the hydra dir, ignored otherwise
	9	+
	10	+ # These settings are based on the cheetah run yaml file
	11	+
	12	+ # General configuration overrides
	13	+ trial_length: 250
	14	+ action_noise_std: 0.3
	15	+
	16	+ # Model overrides
	17	+ num_grad_updates: 100
	18	+ sequence_length: 50
	19	+ batch_size: 50
	20	+ free_nats: 3
	21	+ kl_scale: 1.0
	22	+
	23	+ # Planner configuration overrides
	24	+ planning_horizon: 12
	25	+ cem_num_iters: 10
	26	+ cem_elite_ratio: 0.1
	27	+ cem_population_size: 1000
	28	+ cem_alpha: 0.0
	29	+ cem_clipped_normal: true
	30	+
	31	+ #Ali's overrides
	32	+ model_learning_rate: 1e-4

✓ ⌂ 3 mbrl/examples/conf/overrides/planet\_duckietown.yaml ↗

21	21	cem_population_size: 1000
22	22	cem_alpha: 0.0
23	23	cem_clipped_normal: true
	24	+
	25	+ #Ali's overrides

26 + model\_learning\_rate: 1e-4

✓ ⌂ 3 mbrl/examples/main.py □

```
11   11 import mbrl.algorithms.mppo as mppo
12   12 import mbrl.algorithms.pets as pets
13   13 import mbrl.algorithms.planet as planet
14 + 14 + import mbrl.algorithms.dreamer as dreamer #added April 2022 for project
15   15 import mbrl.util.env
16   16
17   17 import pandas as pd
18   18     return mppo.train(env, test_env, term_fn, cfg)
19   19 if cfg.algorithm.name == "planet":
20   20     return planet.train(env, cfg)
21 + 21 + if cfg.algorithm.name == "dreamer": #added for project
22 + 22 +     return dreamer.train(env, cfg)
23   23
24   24
25   25 if __name__ == "__main__":
```

✓ ⌂ 1 mbrl/models/\_\_init\_\_.py □

```
9   9 from .model_trainer import ModelTrainer
10 10 from .one_dim_tr_model import OneDTransitionRewardModel
11 11 from .planet import PlaNetModel
12 + 12 + from .dreamer import DreamerModel
13   13 from .util import (
14   14     Conv2dDecoder,
15   15     Conv2dEncoder,
```

✓ 723 mbrl/models/dreamer.py □

```
... ... @@ -0,0 +1,723 @@
1 + #####
2 + ##### FOR GLEN #####
3 + #####
4 + # This file holds the bulk of the Dreamer implementation. Most of the code was
5 + # inspired from two public implementations. The original Dreamer author's tf
6 + # implementation here https://github.com/danijar/dreamer and Chandramouli
7 + # Rajagopalan's pytorch implementation of Dreamer here
8 + # https://github.com/chamorajg/pl-dreamer.
9 + #####
10 + from dataclasses import dataclass
11 + from typing import Any, Dict, List, Optional, Tuple, Union
12 +
13 + import numpy as np
14 + import torch
15 + import torch.distributions
16 + import torch.nn as nn
17 + import torch.nn.functional as F
18 +
19 + from mbrl.types import TensorType, TransitionBatch
20 +
21 + from .model import Model
22 + from .util import Conv2dDecoder, Conv2dEncoder, to_tensor #From mbrl-lib.PlaNet
23 +
24 + from PIL import Image
25 +
26 + import wandb
27 + from collections import Iterable
28 + import omegaconf
29 +
30 + from pathlib import Path
31 + import cv2
32 +
33 + from mbrl.models.planet import PlaNetModel #will need ActionDecoder and DenseModel
34 + #from mbrl.models.action import ActionDecoder
35 + #from mbrl.models.dense import DenseModel
```

```

36 + #https://github.com/chamorajg/pl-dreamer/blob/main/dreamer.py
37 + #https://github.com/juliusfrost/dreamer-pytorch/blob/master/dreamer/models/dense.py
38 +
39 + #For Dreamer implementation, Dreamer trainer uses Pytorch Lightning
40 + from tqdm import tqdm
41 + from typing import Callable, Iterator, Tuple
42 + from torch.optim import Adam
43 + from torch.utils.data import DataLoader, IterableDataset
44 + from torch.distributions import Categorical, Normal
45 + from mbrl.models.planet_imp import PLANet, FreezeParameters
46 + from mbrl.models.planet_legacy import Episode, DMControlSuiteEnv
47 +
48 + def flatten_config(cfg, curr_nested_key):
49 +     """The nested config file provided by Hydra cannot be parsed by wandb. This recursive function flattens the
50 +
51 +     Args:
52 +         cfg (Hydra config): The nested config file used by Hydra.
53 +         curr_nested_key (str): The current parent key (used for recursive calls).
54 +
55 +     Returns:
56 +         (dict): A flat configuration dictionary.
57 +     """
58 +
59 +     flat_cfg = {}
60 +
61 +     for curr_key in cfg.keys():
62 +
63 +         # deal with missing values
64 +         try:
65 +             curr_item = cfg[curr_key]
66 +         except Exception as e:
67 +             curr_item = 'NA'
68 +
69 +         # deal with lists
70 +         if type(curr_item) == list or type(curr_item) == omegaconf.listconfig.ListConfig:
71 +             for nested_idx, nested_item in enumerate(curr_item):
72 +                 list_nested_key = f'{curr_nested_key}>{curr_key}>{nested_idx}'
73 +                 flat_cfg[list_nested_key] = nested_item
74 +
75 +         # check if item is also a config
76 +         # recurse
77 +         elif isinstance(curr_item, Iterable) and type(curr_item) != str:
78 +             flat_cfg.update(flatten_config(curr_item, f'{curr_nested_key}>{curr_key}'))
79 +
80 +         # otherwise just add to return dict
81 +         else:
82 +             flat_cfg[f'{curr_nested_key}>{curr_key}'] = curr_item
83 +
84 +     return flat_cfg
85 +
86 + class ExperienceSourceDataset(IterableDataset):
87 +     """
88 +     Implementation from PyTorch Lightning Bolts:
89 +     https://github.com/PyTorchLightning/pytorch-lightning-bolts/blob/master/pl_bolts/datamodules/experience_sour
90 +     Basic experience source dataset. Takes a generate_batch function that returns an iterator.
91 +     The logic for the experience source and how the batch is generated is defined the Lightning model itself
92 +     """
93 +
94 +     def __init__(self, generate_batch: Callable):
95 +         self.generate_batch = generate_batch
96 +
97 +     def __iter__(self) -> Iterator:
98 +         iterator = self.generate_batch
99 +         return iterator
100 +
101 + # class DreamerModel(nn.Module):
102 + class DreamerModel(Model):
103 +

```

```
104 +     def __init__(
105 +         self,
106 +         obs_shape,
107 +         action_size,
108 +         hidden_size_fcs,
109 +         depth_size,
110 +         stoch_size,
111 +         deter_size,
112 +         device: Union[str, torch.device],
113 +         min_std,
114 +     ):
115 +         # This config is needed for now as we figure out
116 +         # what parameters we need to run dreamer
117 +         # we get values from it from time to time
118 +
119 +         outside_config = {
120 +             'name': 'Dreamer',
121 +             'env': 'quadruped_run',
122 +             'seed': 42,
123 +             'ckpt_callback': {
124 +                 'save_top_k': 2,
125 +                 'mode': 'min',
126 +                 'monitor': 'loss',
127 +                 'save_on_train_epoch_end': True,
128 +                 'save_last': True,
129 +                 'trainer_params': None,
130 +                 'default_root_dir': 'None',
131 +                 'gpus': 1,
132 +                 'gradient_clip_val': 100.0,
133 +                 'val_check_interval': 5,
134 +                 'max_epochs': 1000,
135 +             },
136 +             'dreamer': {
137 +                 'td_model_lr': 0.0005,
138 +                 'actor_lr': 8e-05,
139 +                 'critic_lr': 8e-05,
140 +                 'default_lr': 0.0005,
141 +                 'weight_decay': 1e-06,
142 +                 'batch_size': 50,
143 +                 'batch_length': 50,
144 +                 'length': 50,
145 +                 'prefill_timesteps': 5000,
146 +                 'explore_noise': 0.3,
147 +                 'max_episode_length': 1000,
148 +                 'collect_interval': 100,
149 +                 'max_experience_size': 1000,
150 +                 'save_episodes': False,
151 +                 'discount': 0.99,
152 +                 'lambda': 0.95,
153 +                 'clip_actions': False,
154 +                 'horizon': 1000,
155 +                 'imagine_horizon': 15,
156 +                 'free_nats': 3.0,
157 +                 'kl_coeff': 1.0,
158 +                 'dreamer_model': {
159 +                     'obs_space': [3, 64, 64],
160 +                     'num_outputs': 1,
161 +                     'custom_model': 'DreamerModel',
162 +                     'deter_size': 200,
163 +                     'stoch_size': 30,
164 +                     'depth_size': 32,
165 +                     'hidden_size': 400,
166 +                     'action_init_std': 5.0,
167 +                 },
168 +                 'env_config': {'action_repeat': 2},
169 +             },
170 +         }
171 }
```

```

172 +         super().__init__(device)
173 +         self.outside_config = outside_config
174 +
175 +         final_dreamer_config = {"DC" : outside_config}
176 +         flat_cfg = flatten_config(final_dreamer_config, ">")
177 +         for config_item in flat_cfg:
178 +             wandb.config[config_item] = flat_cfg[config_item]
179 +
180 +             self.obs_shape = obs_shape
181 +             self.hidden_size_fcs = hidden_size_fcs
182 +             self.device = device
183 +             self.num_outputs = 1
184 +
185 +             self.model_config = {
186 +                 "hidden_size": self.hidden_size_fcs,
187 +                 'deter_size': deter_size,
188 +                 'stoch_size': stoch_size,
189 +                 'depth_size': depth_size,
190 +                 'action_init_std': min_std,
191 +             }
192 +             self.name = 'Dreamer'
193 +
194 +
195 +     def setGymEnv(self, env, workdir):
196 +         self.env = env
197 +         sample_action_space = np.zeros(self.env.action_space.shape)
198 +         # In the future we may want to use the MBRL-Lib
199 +         # implementation of Planet here
200 +
201 +         # self.model = PlaNetModel #try this when we get the functions mapped properly
202 +
203 +         self.model = PLANet(obs_space= self.obs_shape,
204 +                             action_space= sample_action_space,
205 +                             num_outputs= self.num_outputs,
206 +                             model_config= self.model_config,
207 +                             name = self.name,
208 +                             device = self.device
209 +                         )
210 +         self.episodes = []
211 +         self.length = self.outside_config["dreamer"]['length']
212 +         self.timesteps = 0
213 +         self._max_experience_size = self.outside_config["dreamer"]['max_experience_size']
214 +         self._action_repeat = self.outside_config["dreamer"]["env_config"]["action_repeat"]
215 +         self._prefill_timesteps = self.outside_config["dreamer"]["prefill_timesteps"]
216 +         self._max_episode_length = self.outside_config["dreamer"]["max_episode_length"]
217 +
218 +         self.explore = self.outside_config["dreamer"]['explore_noise']
219 +         self.batch_size = self.outside_config["dreamer"]["batch_size"]
220 +         self.action_space = sample_action_space.shape[0]
221 +         self.imagine_horizon = self.outside_config['dreamer']["imagine_horizon"]
222 +         prefill_episodes = self._prefill_train_batch()
223 +         self._add(prefill_episodes)
224 +         self.workdir = workdir
225 +         self.curr_episode = 0
226 +         self.currently_testing = False
227 +         self.video_counter = 0
228 +         self.in_duckietown = False
229 +
230 +
231 +
232 +     # functions we added to make this work with MBRL-Lib
233 +     def set_curr_episode(self, episode):
234 +         self.curr_episode = episode
235 +
236 +     def set_currently_testing(self, currently_testing):
237 +         self.currently_testing = currently_testing
238

```

```

239 +     def reset_world_model(self, device=None):
240 +         self.model.get_initial_state(device=device)
241 +
242 +     def _process_pixel_obs(self, obs: torch.Tensor) -> torch.Tensor:
243 +         return to_tensor(obs).float().to(self.device) / 256.0 - 0.5
244 +
245 +     def _process_batch(
246 +         self, batch: TransitionBatch, as_float: bool = True, pixel_obs: bool = False
247 +     ) -> Tuple[torch.Tensor, ...]:
248 +         # `obs` is a sequence, so `next_obs` is not necessary
249 +         # sequence iterator samples full sequences, so `dones` not necessary either
250 +         obs, action, _, rewards, _ = super().__process_batch(batch, as_float=as_float)
251 +         if pixel_obs:
252 +             obs = self._process_pixel_obs(obs)
253 +         return obs, action, rewards
254 +
255 +     def compute_dreamer_loss(
256 +         obs,
257 +         action,
258 +         reward,
259 +         imagine_horizon,
260 +         discount=0.99,
261 +         lambda_=0.95,
262 +         kl_coeff=1.0,
263 +         free_nats=3.0,
264 +         log=True):
265 +         """Constructs loss for the Dreamer objective
266 +        Args:
267 +            obs (TensorType): Observations (o_t)
268 +            action (TensorType): Actions (a_(t-1))
269 +            reward (TensorType): Rewards (r_(t-1))
270 +            model (TorchModelV2): DreamerModel, encompassing all other models
271 +            imagine_horizon (int): Imagine horizon for actor and critic loss
272 +            discount (float): Discount
273 +            lambda_ (float): Lambda, like in GAE
274 +            kl_coeff (float): KL Coefficient for Divergence loss in model loss
275 +            free_nats (float): Threshold for minimum divergence in model loss
276 +            log (bool): If log, generate gifs
277 +        """
278 +         encoder_weights = list(self.model.encoder.parameters())
279 +         decoder_weights = list(self.model.decoder.parameters())
280 +         reward_weights = list(self.model.reward.parameters())
281 +         dynamics_weights = list(self.model.dynamics.parameters())
282 +         critic_weights = list(self.model.value.parameters())
283 +         model_weights = list(encoder_weights + decoder_weights + reward_weights +
284 +                               dynamics_weights)
285 +
286 +         device = self.device
287 +         # PlanET Model Loss
288 +         latent = self.model.encoder(obs)
289 +         istate = self.model.dynamics.get_initial_state(obs.shape[0], self.device)
290 +         post, prior = self.model.dynamics.observe(latent, action, istate)
291 +         features = self.model.dynamics.get_feature(post)
292 +         image_pred = self.model.decoder(features)
293 +         reward_pred = self.model.reward(features)
294 +         image_loss = -torch.mean(image_pred.log_prob(obs))
295 +         reward_loss = -torch.mean(reward_pred.log_prob(reward))
296 +         prior_dist = self.model.dynamics.get_dist(prior[0], prior[1])
297 +         post_dist = self.model.dynamics.get_dist(post[0], post[1])
298 +         div = torch.mean(
299 +             torch.distributions.kl_divergence(post_dist, prior_dist).sum(dim=2))
300 +         div = torch.clamp(div, min=free_nats)
301 +         model_loss = kl_coeff * div + reward_loss + image_loss
302 +
303 +         # Actor Loss
304 +         with torch.no_grad():
305

```

```

306 +         actor_states = [v.detach() for v in post]
307 +     with FreezeParameters(model_weights):
308 +         imag_feat = self.model.imagine_ahead(actor_states, imagine_horizon)
309 +     with FreezeParameters(model_weights + critic_weights):
310 +         reward = self.model.reward(imag_feat).mean
311 +         value = self.model.value(imag_feat).mean
312 +     pcont = discount * torch.ones_like(reward)
313 +     returns = self._lambda_return(reward[:-1], value[:-1], pcont[:-1], value[-1],
314 +                                     lambda_)
315 +     discount_shape = pcont[:1].size()
316 +     discount = torch.cumprod(
317 +         torch.cat([torch.ones(*discount_shape).to(device), pcont[:-2]]), dim=0,
318 +         dim=0)
319 +
320 +     # Critic Loss
321 +     with torch.no_grad():
322 +         val_feat = imag_feat.detach()[:-1]
323 +         target = returns.detach()
324 +         val_discount = discount.detach()
325 +         val_pred = self.model.value(val_feat)
326 +         critic_loss = -torch.mean(val_discount * val_pred.log_prob(target))
327 +
328 +     # Logging purposes
329 +     prior_ent = torch.mean(prior_dist.entropy())
330 +     post_ent = torch.mean(post_dist.entropy())
331 +
332 +     log_gif = None
333 +
334 +     if log:
335 +         log_gif = self._log_summary(obs, action, latent, image_pred)
336 +
337 +     return_dict = {
338 +         "model_loss": model_loss,
339 +         "reward_loss": reward_loss,
340 +         "image_loss": image_loss,
341 +         "divergence": div,
342 +         "actor_loss": actor_loss,
343 +         "critic_loss": critic_loss,
344 +         "prior_ent": prior_ent,
345 +         "post_ent": post_ent,
346 +     }
347 +
348 +     if log_gif is not None:
349 +         return_dict["log_gif"] = self._postprocess_gif(log_gif)
350 +     return return_dict
351 +
352 +     # Loss function for dreamer
353 +     # Different from Planet cause more networks
354 +     def loss(
355 +         self,
356 +         batch: TransitionBatch,
357 +         target: Optional[torch.Tensor] = None,
358 +         reduce: bool = True,
359 +     ) -> Tuple[torch.Tensor, Dict[str, Any]]:
360 +         """Computes the Dreamer loss given a batch of transitions.
361 +
362 +         """
363 +         obs, action, rewards = self._process_batch(batch, pixel_obs=True)
364 +
365 +         # hacky check to see if we are using the duckietown environment
366 +         # I just see if the obs shape ends in 3, for the other envs it does not
367 +         if obs.shape[-1] == 3:
368 +             self.in_duckietown = True
369 +
370 +             if self.in_duckietown:
371 +                 obs = obs.permute((0, 1, 4, 2, 3))
372 +

```

```

373 +         return_dict = self.compute_dreamer_loss(obs, action, rewards, self.imagine_horizon)
375 +
376 +         dreamer_obs_loss = return_dict["image_loss"]
377 +         dreamer_reward_loss = return_dict["reward_loss"]
378 +         dreamer_kl_loss = return_dict["divergence"]
379 +
380 +         meta = {
381 +             "reconstruction": None,
382 +             "observations_loss": dreamer_obs_loss.detach().mean().item(),
383 +             "reward_loss": dreamer_reward_loss.detach().mean().item(),
384 +             "kl_loss": dreamer_kl_loss.detach().mean().item(),
385 +         }
386 +
387 +         return return_dict["model_loss"], meta
388 +
389 +     def eval_score(
390 +         self, batch: TransitionBatch, target: Optional[torch.Tensor] = None
391 +     ) -> Tuple[torch.Tensor, Dict[str, Any]]:
392 +         """Computes an evaluation score for the model over the given input/target.
393 +
394 +         This is equivalent to calling loss(batch, reduce=False).
395 +         """
396 +         with torch.no_grad():
397 +             return self.loss(batch, reduce=False)
398 +
399 +         # again, Dreamer update is fundamentally different from
400 +         # from planet update, which causes problems when trying
401 +         # to integrate it into this library
402 +     def dreamer_update(self, dreamer_loss):
403 +
404 +         self.dreamer_optim.zero_grad()
405 +         dreamer_loss.backward()
406 +         self.dreamer_optim.step()
407 +
408 +         return dreamer_loss
409 +
410 +     def dreamer_loss(self, train_batch):
411 +         """ calculates dreamer loss."""
412 +
413 +         log_gif = False
414 +         if "log_gif" in train_batch:
415 +             log_gif = True
416 +
417 +         self.stats_dict = self.compute_dreamer_loss(
418 +             train_batch["obs"],
419 +             train_batch["actions"],
420 +             train_batch["rewards"],
421 +             self.outside_config["dreamer"]["imagine_horizon"],
422 +             self.outside_config["dreamer"]["discount"],
423 +             self.outside_config["dreamer"]["lambda"],
424 +             self.outside_config["dreamer"]["kl_coeff"],
425 +             self.outside_config["dreamer"]["free_nats"],
426 +             log_gif,
427 +         )
428 +
429 +         loss_dict = self.stats_dict
430 +         return loss_dict
431 +
432 +     def _prefill_train_batch(self, ):
433 +         """ Prefill episodes before the training begins."""
434 +
435 +         self.timesteps = 2
436 +         obs = self.env.reset()
437 +         episode = Episode(obs, self.action_space)
438 +         episodes = []
439 +
440 +         while self.timesteps < self._prefill_timesteps:

```

```

441 +         action, logp, state = self._prefill_action_sampler_fn(None,
442 +                                         self.timesteps)
443 +
444 +         action = action.squeeze()
445 +         obs, reward, done, _ = self.env.step(action.numpy())
446 +         episode.append((obs, action, reward, done))
447 +         self.timesteps += self._action_repeat
448 +         if done or self.timesteps == self._prefill_timesteps - 1:
449 +             episodes.append(episode.todict())
450 +             obs = self.env.reset()
451 +             if done:
452 +                 episode.reset(obs)
453 +             del episode
454 +             return episodes
455 +
456 +     def _data_collect(self):
457 +         """ Collect data from the policy after every epoch. """
458 +
459 +         obs = self.env.reset()
460 +         state = self.model.get_initial_state(self.device)
461 +         episode = Episode(obs, self.action_space)
462 +         episodes = []
463 +
464 +         max_len = self._max_episode_length // self._action_repeat
465 +         for i in range(max_len):
466 +             action, logp, state = self._action_sampler_fn(
467 +                 ((episode.obs[-1] / 255.0) - 0.5).unsqueeze(0).to(
468 +                     self.device), state, self.explore, False)
469 +             obs, reward, done, _ = self.env.step(action.detach().cpu().numpy())
470 +             episode.append((obs, action.detach().cpu(), reward, done))
471 +             if done or i == max_len - 1:
472 +                 episodes.append(episode.todict())
473 +                 break
474 +             del episode
475 +         return episodes
476 +
477 +     def _test(self):
478 +         """ Test the model after every few intervals."""
479 +
480 +         obs = self.env.reset()
481 +         state = self.model.get_initial_state(self.device)
482 +         obs = torch.FloatTensor(np.ascontiguousarray(obs.transpose((2, 0, 1))))
483 +
484 +         tot_reward = 0
485 +         done = False
486 +         while not done:
487 +             action, logp, state = self._action_sampler_fn(
488 +                 ((obs / 255.0) - 0.5).unsqueeze(0).to(self.device), state, self.explore, True)
489 +             obs, reward, done, _ = self.env.step(action.detach().cpu().numpy())
490 +             obs = obs.transpose((2, 0, 1))
491 +             obs = torch.FloatTensor(np.ascontiguousarray(obs))
492 +             tot_reward += reward
493 +         return tot_reward
494 +
495 +     def _add(self, batch):
496 +         """ Adds the collected episode samples as well as the pre-filled
497 +             episode samples into the episode memory."""
498 +
499 +         self.episodes.extend(batch)
500 +
501 +         if len(self.episodes) > self._max_experience_size:
502 +             remove_episode_index = len(self.episodes) - \
503 +                 self._max_experience_size
504 +
505 +             if self.outside_config["dreamer"]["save_episodes"] and\
506 +                 self.trainer is not None and self.trainer.log_dir is not None:
507 +                 save_episodes = np.array(self.episodes)
508

```

```

509 +         if not os.path.exists(f'{self.trainer.log_dir}/episodes'):
510 +             os.makedirs(f'{self.trainer.log_dir}/episodes', exist_ok=True)
511 +             np.savez(f'{self.trainer.log_dir}/episodes/episodes.npz', save_episodes)
512 +
513 +     def _sample(self, batch_size):
514 +         """ Samples a batch of episode of length T from the config."""
515 +
516 +         episodes_buffer = []
517 +         while len(episodes_buffer) < batch_size:
518 +             rand_index = random.randint(0, len(self.episodes) - 1)
519 +             episode = self.episodes[rand_index]
520 +             if episode["count"] < self.length:
521 +                 continue
522 +             available = episode["count"] - self.length
523 +             index = int(random.randint(0, available))
524 +             episodes_buffer.append({"count": self.length,
525 +                                     "obs": episode["obs"][index : index + self.length],
526 +                                     "action": episode["action"][index: index + self.length],
527 +                                     "reward": episode["reward"][index: index + self.length],
528 +                                     "done": episode["done"][index: index + self.length],
529 +                                     })
530 +         total_batch = {}
531 +         for k in episodes_buffer[0].keys():
532 +             if k == "count" or k == "state":
533 +                 continue
534 +             else:
535 +                 total_batch[k] = torch.stack([e[k] for e in episodes_buffer], axis=0)
536 +
537 +     def _train_batch(self, batch_size):
538 +         for _ in range(self.outside_config["dreamer"]["collect_interval"]):
539 +             total_batch = self._sample(batch_size)
540 +             def return_batch(i):
541 +                 return (total_batch["obs"][i] / 255.0 - 0.5), \
542 +                        total_batch["action"][i], total_batch["reward"][i], total_batch["done"][i]
543 +                 for i in range(batch_size):
544 +                     yield return_batch(i)
545 +
546 +     def prefill_action_sampler_fn(self, state, timestep):
547 +         """Action sampler function during prefill phase where
548 +         actions are sampled uniformly [-1, 1].
549 +         """
550 +         # Custom Exploration
551 +         logp = [0.0]
552 +         # Random action in space [-1.0, 1.0]
553 +         action = torch.FloatTensor(1, self.model.action_size).uniform_(-1.0,
554 +                           1.0)
555 +         state = self.model.get_initial_state(self.device)
556 +         return action, logp, state
557 +
558 +     def action_sampler_fn(self, obs, state, explore, test=False):
559 +         """Action sampler during training phase, actions
560 +         are evaluated through DreamerPolicy and
561 +         an additive gaussian is added
562 +         to incentivize exploration."""
563 +
564 +         action, logp, state_new = self.model.policy(obs, state,
565 +                                         explore=not(test))
566 +         if not test:
567 +             action = Normal(action, explore).sample()
568 +             action = torch.clamp(action, min=-1.0, max=1.0)
569 +         return action, logp, state_new
570 +
571 +     def training_step(self, batch, batch_idx):
572 +         """ Trains the model on the samples collected."""
573 +
574 +         obs, action, reward, __ = batch
575

```

```

576 +         loss = self.dreamer_loss({"obs":obs,
577 +                               "actions":action, "rewards":reward,
578 +                               "log_gif": True})
579 +     outputs = []
580 +     for k, v in loss.items():
581 +         if "loss" in k:
582 +             self.log(k, v)
583 +         if k in ["model_loss", "critic_loss", "actor_loss"]:
584 +             outputs.append(v)
585 +     return sum(outputs)
586 +
587 + def training_epoch_end(self, outputs):
588 +     """ Collects data samples after every epoch end and tests the
589 +     model on the environment of maximum length from the config every
590 +     few intervals."""
591 +
592 +     total_loss = 0
593 +     for out in outputs:
594 +         total_loss += out['loss'].item()
595 +     if len(outputs) != 0:
596 +         total_loss /= len(outputs)
597 +     self.log('loss', total_loss)
598 +
599 +     with torch.no_grad():
600 +         data_collection_episodes = self._data_collect()
601 +         self._add(data_collection_episodes)
602 +         data_dict = data_collection_episodes[0]
603 +         self.log('avg_reward_collection', torch.mean(data_dict['reward']))
604 +
605 +     if self.current_epoch > 0 and \
606 +         self.current_epoch % self.outside_config["trainer_params"]["val_check_interval"] == 0:
607 +         self.model.eval()
608 +         episode_reward = self._test()
609 +         self.log('avg_reward_test', episode_reward)
610 +         self.model.train()
611 +
612 +     def _collate_fn(self, batch):
613 +         return_batch = {}
614 +         for k in batch[0].keys():
615 +             if k == 'count':
616 +                 return_batch[k] = torch.LongTensor([data[k] for data in batch])
617 +             else:
618 +                 return_batch[k] = torch.stack([data[k] for data in batch])
619 +         return return_batch
620 +
621 +     def train_dataloader(self) -> DataLoader:
622 +         """Get train loader"""
623 +         dataset = ExperienceSourceDataset(self._train_batch(self.batch_size))
624 +         dataloader = DataLoader(dataset=dataset,
625 +                                batch_size=self.batch_size,
626 +                                pin_memory=True,
627 +                                num_workers=1)
628 +
629 +     def configure_optimizers(self,):
630 +         """ Configure optimizers."""
631 +
632 +         encoder_weights = list(self.model.encoder.parameters())
633 +         decoder_weights = list(self.model.decoder.parameters())
634 +         reward_weights = list(self.model.reward.parameters())
635 +         dynamics_weights = list(self.model.dynamics.parameters())
636 +         actor_weights = list(self.model.actor.parameters())
637 +         critic_weights = list(self.model.value.parameters())
638 +
639 +         model_opt = Adam(
640 +             [
641 +                 {'params': encoder_weights + decoder_weights + reward_weights + dynamics_weights,
642 +                  'lr':self.outside_config["dreamer"]["td_model_lr"]}],
```

```

643 +         {'params':actor_weights, 'lr':self.outside_config["dreamer"]["actor_lr"]},
644 +         {'params':critic_weights, 'lr':self.outside_config["dreamer"]["critic_lr"]}],
645 +         lr=self.outside_config["dreamer"]["default_lr"],
646 +         weight_decay=self.outside_config["dreamer"]["weight_decay"])
647 +     self.dreamer_optim = model_opt
648 +
649 +     def _postprocess_gif(self, gif: np.ndarray):
650 +         gif = gif.detach().cpu().numpy()
651 +         gif = np.clip(255*gif, 0, 255).astype(np.uint8)
652 +         B, T, C, H, W = gif.shape
653 +         frames = gif.transpose((1, 2, 3, 0, 4)).reshape((1, T, C, H, B * W))
654 +         frames = frames.squeeze(0)
655 +
656 +     def display_image(frame):
657 +         frame = frame.transpose((1, 2, 0))
658 +         return Image.fromarray(frame)
659 +
660 +     # create destination path for movies
661 +     movie_save_folder = f'{self.workdir}/movies'
662 +     Path(movie_save_folder).mkdir(parents=True, exist_ok=True)
663 +     video_name = f"movie_{self.curr_episode}_{int(self.currently_testing)}_{self.video_counter}"
664 +
665 +     self.video_counter += 1
666 +     # # create a videowriter
667 +     # videodims = (frames.shape[-2], frames.shape[-1])
668 +     # fourcc = cv2.VideoWriter_fourcc(*'MPEG')
669 +     # video = cv2.VideoWriter(f"{movie_save_folder}/{video_name}.mp4", fourcc, 10, videodims)
670 +
671 +     # # write to video
672 +     # for frame in list(frames):
673 +     #     curr_img_frame = display_image(frame)
674 +     #     video.write(cv2.cvtColor(np.array(curr_img_frame), cv2.COLOR_RGB2BGR))
675 +     # video.release()
676 +
677 +     if self.video_counter<200 or self.video_counter%50:
678 +         # also save gif
679 +         img, *imgs = [display_image(frame) for frame in list(frames)]
680 +         # img.save(f'{self.trainer.log_dir}/movies/movie_{self.current_epoch}.gif', format='GIF', append_images=imgs,
681 +         # save_all=True, loop=0)
682 +
683 +         img_save_loc = f'{movie_save_folder}/{video_name}.gif'
684 +         img.save(img_save_loc, format='GIF', append_images=imgs,
685 +         save_all=True, loop=0)
686 +
687 +         image_array = [display_image(frame) for frame in list(frames)]
688 +         # images = wandb.Image(image_array, caption=f"{video_name}")
689 +         wandb_saved_gif = wandb.Image(img_save_loc)
690 +         wandb.log({"video": wandb_saved_gif, "global_episode": self.curr_episode})
691 +
692 +     return frames
693 +
694 +     def _log_summary(self, obs, action, embed, image_pred):
695 +         truth = obs[:6] + 0.5
696 +         recon = image_pred.mean[:6]
697 +         istate = self.model.dynamics.get_initial_state(6, self.device)
698 +         init, _ = self.model.dynamics.observe(embed[:6, :5],
699 +                                             action[:6, :5], istate)
700 +         init = [itm[:6, -1] for itm in init]
701 +         prior = self.model.dynamics.imagine(action[:6, 5:], init)
702 +         openl = self.model.decoder(self.model.dynamics.get_feature(prior)).mean
703 +
704 +         mod = torch.cat([recon[:, :5] + 0.5, openl + 0.5], 1)
705 +         error = (mod - truth + 1.0) / 2.0
706 +         return torch.cat([truth, mod, error], 3)
707 +
708 +     def _lambda_return(self, reward, value, pcont, bootstrap, lambda_):
709

```

```

710 +         def agg_fn(x, y):
711 +             return y[0] + y[1] * lambda_ * x
712 +
713 +             next_values = torch.cat([value[1:], bootstrap[None]], dim=0)
714 +             inputs = reward + pcont * next_values * (1 - lambda_)
715 +
716 +             last = bootstrap
717 +             returns = []
718 +             for i in reversed(range(len(inputs))):
719 +                 last = agg_fn(last, [inputs[i], pcont[i]])
720 +                 returns.append(last)
721 +
722 +             returns = list(reversed(returns))
723 +             returns = torch.stack(returns, dim=0)
724 +
725 +             return returns ⊖

```

mbrl/models/model.py

```

124 124 |     from strings to objects with metadata computed by the model
125 125 |             (e.g., reconstructions, entropy, etc.) that will be used for logging.
126 126 |
127 127 |     # Changes to model class because Dreamer is different
128 128 |     def dreamer_update(self):
129 129 |         return None
130 130 |
131 131 |     def update(
132 132 |         self,
133 133 |             (dict): any additional metadata dictionary computed by :meth:`loss`.
134 134 |             """
135 135 |             self.train()
136 136 |             optimizer.zero_grad()
137 137 |             # Changes to model class because Dreamer is different
138 138 |             # optimizer.zero_grad()
139 139 |             loss, meta = self.loss(model_in, target)
140 140 |             loss.backward()
141 141 |             # loss.backward()
142 142 |             if meta is not None:
143 143 |                 with torch.no_grad():
144 144 |                     grad_norm = 0.0
145 145 |                     for p in list(filter(lambda p: p.grad is not None, self.parameters())):
146 146 |                         grad_norm += p.grad.data.norm(2).item() ** 2
147 147 |                     meta["grad_norm"] = grad_norm
148 148 |             optimizer.step()
149 149 |             # optimizer.step()
150 150 |             # Changes to model class because Dreamer is different
151 151 |             loss = self.dreamer_update(loss)
152 152 |
153 153 |             return loss.item(), meta
154 154 |
155 155 |     def reset():

```

mbrl/models/planet\_imp.py

```

... ... @@ -0,0 +1,559 @@
1 + #####
2 + ##### FOR GLEN #####
3 + #####
4 + # This implementation of PlaNet servers as the world model for our Dreamer
5 + # implementation.
6 + #####
7 + import torch
8 + import torch.nn as nn
9 + import torch.nn.functional as F
10 +
11 + #####
12 + ##### FOR GLEN #####
13 + #####
14 + # This is the PlaNet implementation we ported from Chandramouli Rajagopalan's

```

```

15 + # pytorch implementation of Dreamer here https://github.com/chamorajg/pl-dreamer
16 + # which has helper functions for the non-MBRL-Lib PlaNet, and all of the networks
17 + # (Action and Value Nets as well) and are siloed from the dependencies for QA
18 + # of this version of PlaNet outside of MBRL-Lib PlaNet.
19 + #####
20 + import torch.distributions as td
21 +
22 + import numpy as np
23 +
24 + from typing import Any, List, Tuple, Optional
25 + TensorType = Any
26 +
27 + def atanh(x):
28 +     return 0.5 * torch.log((1 + x) / (1 - x))
29 +
30 + class TanhBijector(torch.distributions.Transform):
31 +     def __init__(self):
32 +         super().__init__()
33 +         self.bijective = True
34 +         self.domain = torch.distributions.constraints.real
35 +         self.codomain = torch.distributions.constraints.interval(-1.0, 1.0)
36 +
37 +     @property
38 +     def sign(self): return 1.
39 +
40 +     def _call(self, x): return torch.tanh(x)
41 +
42 +     def _inverse(self, y: torch.Tensor):
43 +         y = torch.where(
44 +             (torch.abs(y) <= 1.),
45 +             torch.clamp(y, -0.99999997, 0.99999997),
46 +             y)
47 +         y = atanh(y)
48 +         return y
49 +
50 +     def log_abs_det_jacobian(self, x, y):
51 +         return 2. * (np.log(2) - x - F.softplus(-2. * x))
52 +
53 + class Reshape(nn.Module):
54 +
55 +     def __init__(self, shape: List):
56 +         super().__init__()
57 +         self.shape = shape
58 +
59 +     def forward(self, x):
60 +         return x.view(*self.shape)
61 +
62 + class Encoder(nn.Module):
63 +     """ As mentioned in the paper, VAE is used
64 +         to calculate the state posterior needed
65 +         for parameter learning in RSSM. The training objective
66 +         here is to create bound on the data. Here losses
67 +         are written only for observations as rewards losses
68 +         follow them. """
69 +
70 +     def __init__(self,
71 +                  device,
72 +                  depth : int = 32,
73 +                  input_channels : Optional[int] = 3
74 +                  ):
75 +         super(Encoder, self).__init__()
76 +         """
77 +             Initialize the parameters of the Encoder.
78 +             Args
79 +                 depth (int) : Number of channels in the first convolution layer.
80 +                 input_channels (int) : Number of channels in the input observation.
81 +             """
82

```

```

+     self.depth = depth
83 +     self.input_channels = input_channels
84 +     self.encoder = nn.Sequential(
85 +         nn.Conv2d(self.input_channels, self.depth, 4, stride=2),
86 +         nn.ReLU(),
87 +         nn.Conv2d(self.depth, self.depth * 2, 4, stride=2),
88 +         nn.ReLU(),
89 +         nn.Conv2d(self.depth * 2, self.depth * 4, 4, stride=2),
90 +         nn.ReLU(),
91 +         nn.Conv2d(self.depth * 4, self.depth * 8, 4, stride=2),
92 +         nn.ReLU(),
93 +     )
94 +     self.encoder = self.encoder.to(device)
95 +
96 + def forward(self, x):
97 +     """ Flatten the input observation [batch, horizon, 3, 64, 64]
98 +         into shape [batch * horizon, 3, 64, 64] before feeding it
99 +         to the input. """
100 +    orig_shape = x.shape
101 +    x = x.reshape(-1, *x.shape[-3:])
102 +    x = self.encoder(x)
103 +    x = x.reshape(*orig_shape[:-3], -1)
104 +    return x
105 +
106 +
107 + class Decoder(nn.Module):
108 +     """
109 +     Takes the input from the RSSM model
110 +     and then decodes it back to images from
111 +     the latent space model. It is mainly used
112 +     in calculating losses.
113 +     """
114 +     def __init__(self, device,
115 +                  input_size : int,
116 +                  depth: int = 32,
117 +                  shape: Tuple[int] = (3, 64, 64)):
118 +         super(Decoder, self).__init__()
119 +         self.depth = depth
120 +         self.shape = shape
121 +         self.decoder = nn.Sequential(
122 +             nn.Linear(input_size, 32 * self.depth),
123 +             Reshape([-1, 32 * self.depth, 1, 1]),
124 +             nn.ConvTranspose2d(32 * self.depth, 4 * self.depth, 5, stride=2),
125 +             nn.ReLU(),
126 +             nn.ConvTranspose2d(4 * self.depth, 2 * self.depth, 5, stride=2),
127 +             nn.ReLU(),
128 +             nn.ConvTranspose2d(2 * self.depth, self.depth, 6, stride=2),
129 +             nn.ReLU(),
130 +             nn.ConvTranspose2d(self.depth, self.shape[0], 6, stride=2),
131 +         )
132 +         self.decoder = self.decoder.to(device)
133 +
134 +     def forward(self, x):
135 +         orig_shape = x.shape
136 +         x = self.decoder(x)
137 +         reshape_size = orig_shape[:-1] + self.shape
138 +         mean = x.view(*reshape_size)
139 +         return td.Independent( td.Normal(mean, 1), len(self.shape))
140 +
141 +
142 + class ActionDecoder(nn.Module):
143 +     """
144 +     ActionDecoder is the policy module in Dreamer.
145 +
146 +     It outputs a distribution parameterized by mean and std, later to be
147 +     transformed by a custom TanhBijection.
148 +     """
149

```

```

150 +     def __init__(self,
151 +         device,
152 +         input_size: int,
153 +         action_size: int,
154 +         layers: int,
155 +         units: int,
156 +         dist: str = "tanh_normal",
157 +         min_std: float = 1e-4,
158 +         init_std: float = 5.0,
159 +         mean_scale: float = 5.0):
160 +     super(ActionDecoder, self).__init__()
161 +     self.layers = layers
162 +     self.units = units
163 +     self.dist = dist
164 +     self.act = nn.ReLU
165 +     self.min_std = min_std
166 +     self.init_std = init_std
167 +     self.mean_scale = mean_scale
168 +     self.action_size = action_size
169 +
170 +     self.layers = []
171 +     self.softplus = nn.Softplus()
172 +
173 +     # MLP Construction
174 +     cur_size = input_size
175 +     for _ in range(self.layers):
176 +         self.layers.extend([nn.Linear(cur_size, self.units), self.act()])
177 +         cur_size = self.units
178 +     self.model = nn.Sequential(*self.layers)
179 +     self.model = self.model.to(device)
180 +
181 +     def forward(self, x):
182 +         raw_init_std = np.log(np.exp(self.init_std) - 1)
183 +         x = self.model(x)
184 +         mean, std = torch.chunk(x, 2, dim=-1)
185 +         mean = self.mean_scale * torch.tanh(mean / self.mean_scale)
186 +         std = self.softplus(std + raw_init_std) + self.min_std
187 +         dist = td.Normal(mean, std)
188 +         transforms = [TanhBijector()]
189 +         dist = td.transformed_distribution.TransformedDistribution(
190 +             dist, transforms)
191 +         dist = td.Independent(dist, 1)
192 +         return dist
193 +
194 +
195 + class DenseDecoder(nn.Module):
196 +     """
197 +     FC network that outputs a distribution for calculating log_prob.
198 +     Used later in DreamerLoss.
199 +     """
200 +
201 +     def __init__(self,
202 +         device,
203 +         input_size: int,
204 +         output_size: int,
205 +         layers: int,
206 +         units: int,
207 +         dist: str = "normal"):
208 +         """Initializes FC network
209 +         Args:
210 +             input_size (int): Input size to network
211 +             output_size (int): Output size to network
212 +             layers (int): Number of layers in network
213 +             units (int): Size of the hidden layers
214 +             dist (str): Output distribution, parameterized by FC output
215 +             logits.
216

```

```

217 +         act (Any): Activation function
218 +
219 +     """
220 +     super().__init__()
221 +     self.layers = layers
222 +     self.units = units
223 +     self.act = nn.ELU
224 +     self.dist = dist
225 +     self.input_size = input_size
226 +     self.output_size = output_size
227 +     self.layers = []
228 +     cur_size = input_size
229 +     for _ in range(self.layers):
230 +         self.layers.extend([nn.Linear(cur_size, self.units), self.act()])
231 +         cur_size = units
232 +     self.layers.append(nn.Linear(cur_size, output_size))
233 +     self.model = nn.Sequential(*self.layers)
234 +     self.model = self.model.to(device)
235 +
236 +
237 +     def forward(self, x):
238 +         x = self.model(x)
239 +         if self.output_size == 1:
240 +             x = torch.squeeze(x)
241 +         if self.dist == "normal":
242 +             output_dist = td.Normal(x, 1)
243 +         elif self.dist == "binary":
244 +             output_dist = td.Bernoulli(logits=x)
245 +         else:
246 +             raise NotImplementedError("Distribution type not implemented!")
247 +         return td.Independent(output_dist, 0)
248 +
249 +
250 +     class RSSM(nn.Module):
251 +         """RSSM is the core recurrent part of the PlaNET module. It consists of
252 +         two networks, one (obs) to calculate posterior beliefs and states and
253 +         the second (img) to calculate prior beliefs and states. The prior network
254 +         takes in the previous state and action, while the posterior network takes
255 +         in the previous state, action, and a latent embedding of the most recent
256 +         observation.
257 +
258 +         Args:
259 +             device,
260 +             action_size: int,
261 +             embed_size: int,
262 +             stoch: int = 30,
263 +             deter: int = 200,
264 +             hidden: int = 200):
265 +             """Initializes RSSM
266 +             Args:
267 +                 action_size (int): Action space size
268 +                 embed_size (int): Size of ConvEncoder embedding
269 +                 stoch (int): Size of the distributional hidden state
270 +                 deter (int): Size of the deterministic hidden state
271 +                 hidden (int): General size of hidden layers
272 +                 act (Any): Activation function
273 +
274 +             super().__init__()
275 +             self.stoch_size = stoch
276 +             self.deter_size = deter
277 +             self.hidden_size = hidden
278 +             self.act = nn.ELU
279 +             # self.act = self.act.to(device)
280 +             self.obs1 = nn.Linear(embed_size + deter, hidden)
281 +             self.obs1 = self.obs1.to(device)
282 +             self.obs2 = nn.Linear(hidden, 2 * stoch)
283 +             self.obs2 = self.obs2.to(device)

```

```

284 +         self.cell = nn.GRUCell(self.hidden_size, hidden_size=self.deter_size)
285 +         self.cell = self.cell.to(device)
286 +         self.img1 = nn.Linear(stoch + action_size, hidden)
287 +         self.img1 = self.img1.to(device)
288 +         self.img2 = nn.Linear(deter, hidden)
289 +         self.img2 = self.img2.to(device)
290 +         self.img3 = nn.Linear(hidden, 2 * stoch)
291 +         self.img3 = self.img3.to(device)
292 +
293 +         self.softplus = nn.Softplus
294 +         # self.softplus = self.softplus.to(device)
295 +
296 +
297 +     def get_initial_state(self, batch_size: int, device) -> List[TensorType]:
298 +
299 +         """Returns the initial state for the RSSM, which consists of mean,
300 +         std for the stochastic state, the sampled stochastic hidden state
301 +         (from mean, std), and the deterministic hidden state, which is
302 +         pushed through the GRUCell.
303 +
304 +         Args:
305 +             batch_size (int): Batch size for initial state
306 +         Returns:
307 +             List of tensors
308 +         """
309 +
310 +         return [
311 +             torch.zeros(batch_size, self.stoch_size).to(device),
312 +             torch.zeros(batch_size, self.stoch_size).to(device),
313 +             torch.zeros(batch_size, self.stoch_size).to(device),
314 +             torch.zeros(batch_size, self.deter_size).to(device),
315 +
316 +
317 +     def observe(self,
318 +                 embed: TensorType,
319 +                 action: TensorType,
320 +                 state: List[TensorType] = None
321 + ) -> Tuple[List[TensorType], List[TensorType]]:
322 +
323 +         """Returns the corresponding states from the embedding from ConvEncoder
324 +         and actions. This is accomplished by rolling out the RNN from the
325 +         starting state through each index of embed and action, saving all
326 +         intermediate states between.
327 +
328 +         Args:
329 +             embed (TensorType): ConvEncoder embedding
330 +             action (TensorType): Actions
331 +             state (List[TensorType]): Initial state before rollout
332 +
333 +         Returns:
334 +             Posterior states and prior states (both List[TensorType])
335 +
336 +         if state is None:
337 +             state = self.get_initial_state(action.size()[0])
338 +
339 +         if embed.dim() <= 2:
340 +             embed = torch.unsqueeze(embed, 1)
341 +
342 +         if action.dim() <= 2:
343 +             action = torch.unsqueeze(action, 1)
344 +
345 +             embed = embed.permute(1, 0, 2)
346 +             action = action.permute(1, 0, 2)
347 +
348 +             priors = [[] for i in range(len(state))]
349 +             posts = [[] for i in range(len(state))]
350 +             last = (state, state)
351 +
352 +             for index in range(len(action)):
353 +                 # Tuple of post and prior
354 +                 last = self.obs_step(last[0], action[index], embed[index])
355 +                 [o.append(s) for s, o in zip(last[0], posts)]
356 +                 [o.append(s) for s, o in zip(last[1], priors)]

```

```

351 +         prior = [torch.stack(x, dim=0) for x in priors]
352 +         post = [torch.stack(x, dim=0) for x in posts]
353 +
354 +         prior = [e.permute(1, 0, 2) for e in prior]
355 +         post = [e.permute(1, 0, 2) for e in post]
356 +
357 +     return post, prior
358 +
359 + def imagine(self, action: TensorType,
360 +             state: List[TensorType] = None) -> List[TensorType]:
361 +     """Imagines the trajectory starting from state through a list of actions.
362 +     Similar to observe(), requires rolling out the RNN for each timestep.
363 +     Args:
364 +         action (TensorType): Actions
365 +         state (List[TensorType]): Starting state before rollout
366 +     Returns:
367 +         Prior states
368 +     """
369 +     if state is None:
370 +         state = self.get_initial_state(action.size()[0])
371 +
372 +     action = action.permute(1, 0, 2)
373 +
374 +     indices = range(len(action))
375 +     priors = [[] for _ in range(len(state))]
376 +     last = state
377 +     for index in indices:
378 +         last = self.img_step(last, action[index])
379 +         [o.append(s) for s, o in zip(last, priors)]
380 +
381 +     prior = [torch.stack(x, dim=0) for x in priors]
382 +     prior = [e.permute(1, 0, 2) for e in prior]
383 +     return prior
384 +
385 + def obs_step(
386 +     self, prev_state: TensorType, prev_action: TensorType,
387 +     embed: TensorType) -> Tuple[List[TensorType], List[TensorType]]:
388 +     """Runs through the posterior model and returns the posterior state
389 +     Args:
390 +         prev_state (TensorType): The previous state
391 +         prev_action (TensorType): The previous action
392 +         embed (TensorType): Embedding from ConvEncoder
393 +     Returns:
394 +         Post and Prior state
395 +     """
396 +     prior = self.img_step(prev_state, prev_action)
397 +     x = torch.cat([prior[3], embed], dim=-1)
398 +     x = self.obs1(x)
399 +     x = self.act()(x)
400 +     x = self.obs2(x)
401 +     mean, std = torch.chunk(x, 2, dim=-1)
402 +     std = self.softplus()(std) + 0.1
403 +     stoch = self.get_dist(mean, std).rsample()
404 +     post = [mean, std, stoch, prior[3]]
405 +     return post, prior
406 +
407 + def img_step(self, prev_state: TensorType,
408 +              prev_action: TensorType) -> List[TensorType]:
409 +     """Runs through the prior model and returns the prior state
410 +     Args:
411 +         prev_state (TensorType): The previous state
412 +         prev_action (TensorType): The previous action
413 +     Returns:
414 +         Prior state
415 +     """
416 +     x = torch.cat([prev_state[2], prev_action], dim=-1)
417

```

```

418 +         x = self.img1(x)
419 +         x = self.act()(x)
420 +         deter = self.cell(x, prev_state[3])
421 +         x = deter
422 +         x = self.img2(x)
423 +         x = self.act()(x)
424 +         x = self.img3(x)
425 +         mean, std = torch.chunk(x, 2, dim=-1)
426 +         std = self.softplus()(std) + 0.1
427 +         stoch = self.get_dist(mean, std).rsample()
428 +         return [mean, std, stoch, deter]
429 +
430 +     def get_feature(self, state: List[TensorType]) -> TensorType:
431 +         # Constructs feature for input to reward, decoder, actor, critic
432 +         return torch.cat([state[2], state[3]], dim=-1)
433 +
434 +     def get_dist(self, mean: TensorType, std: TensorType) -> TensorType:
435 +         return td.Normal(mean, std)
436 +
437 + # Dreamer Model
438 + class PLANet(nn.Module):
439 +     def __init__(self, obs_space, action_space, num_outputs, model_config,
440 +                  name, device):
441 +         super().__init__()
442 +
443 +         nn.Module.__init__(self)
444 +         # 'dreamer_model': {
445 +             # 'obs_space': [3, 64, 64],
446 +             # 'num_outputs': 1,
447 +             # 'custom_model': 'DreamerModel',
448 +             # 'deter_size': 200,
449 +             # 'stoch_size': 30,
450 +             # 'depth_size': 32,
451 +             # 'hidden_size': 400,
452 +             # 'action_init_std': 5.0,
453 +             # },
454 +         self.depth = model_config["depth_size"]
455 +         self.deter_size = model_config["deter_size"]
456 +         self.stoch_size = model_config["stoch_size"]
457 +         self.hidden_size = model_config["hidden_size"]
458 +
459 +         self.action_size = action_space.shape[0]
460 +         self.device = device
461 +
462 +         self.encoder = Encoder(device, self.depth)
463 +         self.decoder = Decoder(device,
464 +                               self.stoch_size + self.deter_size, depth=self.depth)
465 +         self.reward = DenseDecoder(device, self.stoch_size + self.deter_size, 1, 2,
466 +                                   self.hidden_size)
467 +         self.dynamics = RSSM(device,
468 +                               self.action_size,
469 +                               32 * self.depth,
470 +                               stoch=self.stoch_size,
471 +                               deter=self.deter_size)
472 +         self.actor = ActionDecoder(device, self.stoch_size + self.deter_size,
473 +                                   self.action_size, 4, self.hidden_size)
474 +         self.value = DenseDecoder(device, self.stoch_size + self.deter_size, 1, 3,
475 +                                   self.hidden_size)
476 +         self.state = None
477 +
478 +     def policy(self, obs: TensorType, state: List[TensorType], explore=True
479 +               ) -> Tuple[TensorType, List[float], List[TensorType]]:
480 +         """Returns the action. Runs through the encoder, recurrent model,
481 +         and policy to obtain action.
482 +         """
483 +         if state is None:
484

```

```

485 +         # self.state = self.get_initial_state(batch_size=obs.shape[0])
486 +         self.state = self.get_initial_state(self.device)
487 +     else:
488 +         self.state = state
489 +         post = self.state[:4]
490 +         action = self.state[4]
491 +
492 +
493 +         obs = obs.to(self.device)
494 +
495 +         embed = self.encoder(obs)
496 +         post, _ = self.dynamics.obs_step(post, action, embed)
497 +         feat = self.dynamics.get_feature(post)
498 +
499 +         action_dist = self.actor(feat)
500 +         if explore:
501 +             action = action_dist.sample()
502 +         else:
503 +             samples = []
504 +             for _ in range(1000):
505 +                 samples.append(action_dist.sample())
506 +             action = torch.mean(torch.cat(samples), dim=0)
507 +             if action.ndim == 1:
508 +                 action = action.unsqueeze(0)
509 +             logp = action_dist.log_prob(action)
510 +
511 +
512 +     self.state = post + [action]
513 +     return action, logp, self.state
514 +
515 +     def imagine_ahead(self, state: List[TensorType],
516 +                      horizon: int) -> TensorType:
517 +         """Given a batch of states, rolls out more state of length horizon.
518 +         """
519 +         start = []
520 +         for s in state:
521 +             s = s.contiguous().detach()
522 +             shpe = [-1] + list(s.size())[2:]
523 +             start.append(s.view(*shpe))
524 +
525 +         def next_state(state):
526 +             feature = self.dynamics.get_feature(state).detach()
527 +             action = self.actor(feature).rsample()
528 +             next_state = self.dynamics.img_step(state, action)
529 +             return next_state
530 +
531 +         last = start
532 +         outputs = [[] for i in range(len(start))]
533 +         for _ in range(horizon):
534 +             last = next_state(last)
535 +             [o.append(s) for s, o in zip(last, outputs)]
536 +         outputs = [torch.stack(x, dim=0) for x in outputs]
537 +
538 +         imag_feat = self.dynamics.get_feature(outputs)
539 +         return imag_feat
540 +
541 +     def get_initial_state(self, device) -> List[TensorType]:
542 +         self.state = self.dynamics.get_initial_state(1, device) + [
543 +             torch.zeros(1, self.action_size).to(device)
544 +         ]
545 +         return self.state
546 +
547 +
548 +     class FreezeParameters:
549 +         def __init__(self, parameters):
550 +             self.parameters = parameters
551

```

```

552 +         self.param_states = [p.requires_grad for p in self.parameters]
553 +
554 +     def __enter__(self):
555 +         for param in self.parameters:
556 +             param.requires_grad = False
557 +
558 +     def __exit__(self, exc_type, exc_val, exc_tb):
559 +         for i, param in enumerate(self.parameters):
560 +             param.requires_grad = self.param_states[i] ⊖

```

▼ 138 mbrl/models/planet\_legacy.py □

...	...	@@ -0,0 +1,138 @@
-----	-----	-------------------

```

1 + ##### FOR GLEN #####
2 + ##### FOR GLEN #####
3 + #####
4 + # This implementation of PlaNet servers as the world model for our Dreamer
5 + # implementation.
6 +
7 + #####
8 + ##### FOR GLEN #####
9 + #####
10 + # This was used to do QA on the PlaNet implementation we ported from
11 + # Chandramouli Rajagopalan's pytorch implementation of Dreamer here
12 + # https://github.com/chamorajg/pl-dreamer , which has helper functions for the
13 + # non-MBRL-Lib PlaNet, but do not get used with MBRL-Lib or Duckietown.
14 + #####
15 + import torch
16 + import numpy as np
17 + from typing import Tuple, Any, Union, Optional, Dict
18 + import gym
19 + TensorType = Any
20 +
21 + class DMControlSuiteEnv:
22 +
23 +     def __init__(self,
24 +                  name: str,
25 +                  max_episode_length: int = 1000,
26 +                  action_repeat:int = 2,
27 +                  size: Tuple[int] = (64, 64),
28 +                  camera: Optional[Any] = None,
29 +                  ):
30 +         domain, task = name.split('_', 1)
31 +         if domain == 'cup':
32 +             domain = 'ball_in_cup'
33 +         if isinstance(domain, str):
34 +             from dm_control import suite
35 +             self._env = suite.load(domain, task)
36 +         else:
37 +             assert task is None
38 +             self._env = domain()
39 +         self._size = size
40 +         if camera is None:
41 +             camera = dict(quadruped=2).get(domain, 0)
42 +         self._camera = camera
43 +         self._step = 0
44 +         self._max_episode_length = max_episode_length
45 +         self._action_repeat = action_repeat
46 +
47 +     @property
48 +     def observation_space(self):
49 +         spaces = {}
50 +         for key, value in self._env.observation_spec().items():
51 +             spaces[key] = gym.spaces.Box(
52 +                 -np.inf, np.inf, value.shape, dtype=np.float32)
53 +             spaces['image'] = gym.spaces.Box(
54 +                 0, 255, self._size + (3,), dtype=np.uint8)
55 +

```

```

        return gym.spaces.Dict(spaces)

56 +
57 +     @property
58 +     def action_space(self):
59 +         spec = self._env.action_spec()
60 +         return gym.spaces.Box(spec.minimum, spec.maximum, dtype=np.float32)
61 +
62 +     def step(self, action):
63 +         reward = 0
64 +         obs = None
65 +         for k in range(self._action_repeat):
66 +             time_step = self._env.step(action)
67 +             self._step += 1
68 +             obs = dict(time_step.observation)
69 +             obs['image'] = self.render()
70 +             reward += time_step.reward or 0
71 +             done = time_step.last() or self._step == self._max_episode_length
72 +             if done:
73 +                 break
74 +         info = {'discount': np.array(time_step.discount, np.float32)}
75 +         return obs["image"], reward, done, info
76 +
77 +     def reset(self):
78 +         time_step = self._env.reset()
79 +         self._step = 0
80 +         obs = dict(time_step.observation)
81 +         obs['image'] = self.render()
82 +         return obs["image"]
83 +
84 +     def render(self, *args, **kwargs):
85 +         if kwargs.get('mode', 'rgb_array') != 'rgb_array':
86 +             raise ValueError("Only render mode 'rgb_array' is supported.")
87 +         return self._env.physics.render(*self._size, camera_id=self._camera)
88 +
89 + class Episode(object):
90 +     """ Episode Class which contains the related
91 +         attributes of an environment episode in the
92 +         the format similar to queue"""
93 +
94 +     def __init__(self,
95 +                  obs:TensorType,
96 +                  action_space: int = 1,
97 +                  action_repeat: int = 2) -> None:
98 +         """Initializes a list of all episode attributes"""
99 +         self.action_space = action_space
100 +        self.action_repeat = action_repeat
101 +        obs = torch.FloatTensor(np.ascontiguousarray(obs.transpose
102 +                                                    ((2, 0, 1))))
103 +        self.t = 1
104 +        self.obs = [obs]
105 +        self.action = [torch.FloatTensor(torch.zeros(1, self.action_space)).squeeze()]
106 +        self.reward = [0]
107 +        self.done = [False]
108 +
109 +    def append(self,
110 +              episode_attrs: Tuple[TensorType]) -> None:
111 +        """ Appends episode attribute to the list."""
112 +        obs, action, reward, done = episode_attrs
113 +        obs = torch.FloatTensor(np.ascontiguousarray(obs.transpose
114 +                                                    ((2, 0, 1))))
115 +        self.t += 1
116 +        self.obs.append(obs)
117 +        self.action.append(action)
118 +        self.reward.append(reward)
119 +        self.done.append(done)
120 +
121 +    def reset(self,
122 +

```

```

123 +         obs:TensorType) -> None:
124 +         """ Resets Episode list of attributes."""
125 +         obs = torch.FloatTensor(np.ascontiguousarray(obs.transpose
126 +                                         ((2, 0, 1))))
126 +         self.t = 1
127 +         self.obs = [obs]
128 +         self.action = [torch.FloatTensor(torch.zeros(1, self.action_space)).squeeze()]
129 +         self.reward = [0]
130 +         self.done = [False]
131 +
132 +     def todict(self,) -> Dict:
133 +         episode_dict = dict({'count': self.t,
134 +                               'obs': torch.stack(self.obs),
135 +                               'action': torch.cat(self.action),
136 +                               'reward': torch.FloatTensor(self.reward),
137 +                               'done': torch.BoolTensor(self.done)})
138 +
139 +     return ⊖

```

▽ 71 mbrl/planning/dreamer\_wrapper.py □

...	...	@@ -0,0 +1,71 @@
-----	-----	------------------

```

1 + #####
2 + ##### FOR GLEN #####
3 + #####
4 + # This library uses SAC as a "planner" or "agent" for some dynamics models. We
5 + # tried to get Dreamer to be an agent similar to this approach. But the planner
6 + # isn't supposed to have access to the dynamics model for planning, but Dreamer
7 + # needs access to "imagine" future trajectories to determine what action to
8 + # take. We decided it would be simpler to instead put the planning part inside
9 + # the world model, and query the world model for actions as well.
10 +
11 + # We also need to differentiate through the dynamics model, using analytical
12 + # gradients and stop gradients - but were if unsure that the MBRL-Lib accounted
13 + # for this. All of these are design choices needing to be made before being neck
14 + # deep in an implementation, with a commitment made in algorithms/dreamer.py
15 + # before developing models/dreamer.py or planning/dreamer_wrapper.py. It would
16 + # seem usually you take your algorithm and use it in the library, but in this
17 + # case, you need to design your algorithm around MBRL with its training helper
18 + # functions.
19 +
20 + # This was the beginning of our properly MBRL-Lib philosophy- abiding version of
21 + # Dreamer, with a proper agent, based on the SAC implementation shown below.
22 + # Might be possible to add what we did as a third_party extension in this way as
23 + # well.
24 + #####
25 +
26 + # Copyright (c) Facebook, Inc. and its affiliates. All Rights Reserved.
27 + #
28 + # This source code is licensed under the MIT license found in the
29 + # LICENSE file in the root directory of this source tree.
30 + import numpy as np
31 + import torch
32 +
33 + # import mbrl.third_party.pytorch_sac as pytorch_sac
34 + # import mbrl.third_party.pytorch_sac.utils as pytorch_sac_utils
35 +
36 + from .core import Agent
37 +
38 +
39 + class DreamerAgent(Agent):
40 +     """A Soft-Actor Critic agent.
41 +
42 +     This class is a wrapper for
43 +     https://github.com/luisenp/pytorch_sac/blob/master/pytorch_sac/agent/sac.py
44 +
45 +
46 +     Args:
47 +         (pytorch_sac.SACAgent): the agent to wrap.

```

```

48 + """
49 +
50 +     def __init__(self, world_model = None, actor_network = None, critic_network = None):
51 +         self.world_model = world_model
52 +         self.actor_network = actor_network
53 +         self.critic_network = critic_network
54 +
55 +     def act(
56 +         self, obs: np.ndarray, sample: bool = False, batched: bool = False, **_kwargs
57 +     ) -> np.ndarray:
58 +         """Issues an action given an observation.
59 +
60 +         Args:
61 +             obs (np.ndarray): the observation (or batch of observations) for which the action
62 +                 is needed.
63 +             sample (bool): if ``True`` the agent samples actions from its policy, otherwise it
64 +                 returns the mean policy value. Defaults to ``False``.
65 +             batched (bool): if ``True`` signals to the agent that the obs should be interpreted
66 +                 as a batch.
67 +
68 +         Returns:
69 +             (np.ndarray): the action.
70 +         """
71 +         return self.actor_network(obs)

```

```

▼ 15 mbrl/util/env.py □
... ... @@ -1,3 +1,9 @@
1 + ##### FOR GLEN #####
2 + ##### FOR GLEN #####
3 + #####
4 + # So here we just added a fixed start point for the Duckietown environment. That
5 + # required modifying the Duckietown environment as well.
6 + #####
7 # Copyright (c) Facebook, Inc. and its affiliates. All Rights Reserved.
8 #
9 # This source code is licensed under the MIT license found in the
92 # Use the duckietown environment with settings that
93 # closely match the settings used for cheetah run
94 elif cfg.overrides.env == "duckietown_gym_env":
95     env = mbrl.env.DuckietownEnv(domain_rand=False, camera_width=64, camera_height=64)
101 # env = mbrl.env.DuckietownEnv(domain_rand=False, camera_width=64, camera_height=64, set_start_pos=
102 env = mbrl.env.DuckietownEnv(
103     domain_rand=False,
104     camera_width=64,
105     camera_height=64,
106     set_start_pos=np.array([3.68680743, 0.0, 2.01533089]),
107     set_start_angle=np.array([-1.564463624086557])
108 )
96 term_fn = mbrl.env.termination_fns.no_termination
97 reward_fn = None
98 else:

```

```

▼ 1 requirements/ducky.txt □
21 21 torch==1.10.1
22 22 torchvision==0.11.2
23 23 tqdm==4.62.3
24 + typing_extensions==4.1.1
25 wandb==0.12.9

```

```

▼ 30 sweep.py □
... ... @@ -0,0 +1,30 @@
1 + ##### FOR GLEN #####
2 + ##### FOR GLEN #####
3 + #####

```

```
4 + # This is just a script to run hyperparameter sweeps using weights and biases.
5 + #####
6 +
7 + import os
8 + import sys
9 +
10 + base_command = f"python -m mbrl.examples.main algorithm=planet dynamics_model=planet overrides=planet_duckietown"
11 +
12 + final_hydra_command = base_command
13 +
14 + for wandb_arg in sys.argv[1:]:
15 +     hydra_arg = wandb_arg[2:]
16 +     final_hydra_command += ' ' + f"{hydra_arg}"
17 +
18 +
19 + print(final_hydra_command)
20 + os.system(final_hydra_command)
21 +
22 + # Alternate base commands:
23 +
24 + # python -m mbrl.examples.main algorithm=planet dynamics_model=planet overrides=planet_duckietown algorithm.test
25 +
26 + # python -m mbrl.examples.main algorithm=dreamer dynamics_model=dreamer overrides=dreamer_cheetah_run algorithm.
27 +
28 + # python -m mbrl.examples.main algorithm=planet dynamics_model=planet overrides=dreamer_duckietown algorithm.tes
29 +
30 + # python -m mbrl.examples.main algorithm=dreamer dynamics_model=dreamer overrides=dreamer_duckietown algorithm.t
```

## Comparing changes

This is a direct comparison between two commits made in this repository or its related repositories. View the default comparison for this range [here](#).



base: 3af2f04 ▾



compare: 8fa9fc6 ▾

Showing 9 changed files with 262 additions and 12 deletions.

[Split](#) [Unified](#)

▼ 2 .gitignore □

7	7	bin
8	8	
9	9	trained_models
	10	+ saved_runs
	11	+ wandb ⊖

▼ 21 .vscode/launch.json □

...	...	@@ -0,0 +1,21 @@
1		+ {
2		+ // Use IntelliSense to learn about possible attributes.
3		+ // Hover to view descriptions of existing attributes.
4		+ // For more information, visit: https://go.microsoft.com/fwlink/?linkid=830387
5		"version": "0.2.0",
6		"configurations": [
7		{
8		"name": "Manual Control",
9		"type": "python",
10		"request": "launch",
11		"python": "/home/kuwajerw/anaconda3/envs/dgym/bin/python",
12		"program": "manual_control.py",
13		"args": [
14		"--env-name",
15		"Duckietown-udem1-v0"
16		],
17		"console": "integratedTerminal",
18		"justMyCode": true
19		}
20		]
21		+ } ⊖

▼ 34 gym\_duckietown/simulator.py □

...	...	@@ -1,12 +1,14 @@
1	1	# coding=utf-8
2	2	from __future__ import division
3	3	
4		- from collections import namedtuple
	4	+ from collections import namedtuple, deque
5	5	from ctypes import POINTER
6	6	from dataclasses import dataclass
7	7	from typing import Tuple
8	8	# import geometry
9	9	
10		+ # import wandb
11		+
12	12	@dataclass
13	13	class DoneRewardInfo:
14	14	done: bool
158	160	seed=None,

```

159 161         distortion=False,
160 162         randomize_maps_on_reset=False,
163 +     set_start_pos=None,
164 +     set_start_angle=None
161 165     ):
162 166     """
163 167
305 309         self.map_names = os.listdir('maps')
306 310         self.map_names = [mapfile.replace('.yaml', '') for mapfile in self.map_names]
307 311
312 +
313 +     self.set_start_pos = set_start_pos
314 +     self.set_start_angle = set_start_angle
315 +
308 316     # Initialize the state
309 317     self.reset()
310 318
513 521     self.cur_pos = propose_pos
514 522     self.cur_angle = propose_angle
515 523
524 +     if self.set_start_pos is not None:
525 +         self.cur_pos = self.set_start_pos
526 +     if self.set_start_angle is not None:
527 +         # self.cur_angle = self.set_start_angle
528 +         # v Dunno why I need to do this but ugh v
529 +         self.cur_angle = np.array([-1.56446362])
530 +         xyz = 5
531 +
532 +     self.cur_poses = deque(maxlen=10)
533 +     self.cur_angles = deque(maxlen=10)
534 +     self.cur_poses.append(self.cur_pos)
535 +     self.cur_angles.append(self.cur_angle)
536 +
516 537     logger.info('Starting at %s %s' % (self.cur_pos, self.cur_angle))
517 538
518 539     # Generate the first camera image
1220 1241
1221 1242             self.wheel_dist,
1222 1243             wheelVels=self.wheelVels,
                               deltaTime=delta_time)
1244 +
1245 +     self.cur_poses.append(self.cur_pos)
1246 +     self.cur_angles.append(self.cur_angle)
1247 +
1223 1248     self.step_count += 1
1224 1249     self.timestamp += delta_time
1225 1250
1324 1349
1325 1350     # Compute the reward
1326 1351     reward = (
1327 -         +1.0 * speed * lp.dot_dir +
1328 1353             -10 * np.abs(lp.dist) +
1329 1354             +40 * col_penalty
1330 1355         )
1356 +     # Ali's reward:
1357 +     delta_pos = self.cur_pos - self.cur_poses[0]
1358 +     distance_travelled = np.linalg.norm(delta_pos)
1359 +     # wandb.log({"dist_traveled": distance_traveled})
1360 +     reward += distance_travelled*50
1331 1361     return reward
1332 1362
1333 1363     def step(self, action: np.ndarray):

```

v ⌂ 4 learning/reinforcement/pytorch/ddpg.py □

```

67 67     x = self.bn2(self.lr(self.conv2(x)))
68 68     x = self.bn3(self.lr(self.conv3(x)))
69 69     x = self.bn4(self.lr(self.conv4(x)))

```

70	70	-       x = x.view(x.size(0), -1) # flatten +       x = x.reshape(x.size(0), -1) # flatten
71	71	x = self.dropout(x)
72	72	x = self.lr(self.lin1(x))
73	73	
129	129	x = self.bn2(self.lr(self.conv2(x)))
130	130	x = self.bn3(self.lr(self.conv3(x)))
131	131	x = self.bn4(self.lr(self.conv4(x)))
132	132	-       x = x.view(x.size(0), -1) # flatten +       x = x.reshape(x.size(0), -1) # flatten
133	133	x = self.lr(self.lin1(x))
134	134	x = self.lr(self.lin2(torch.cat([x, actions], 1))) # c
135	135	x = self.lin3(x)

▼ ⌂ 6 learning/reinforcement/pytorch/enjoy\_reinforcement.py ↗

6	6	import numpy as np
7	7	
8	8	# Duckietown Specific
9	9	- from reinforcement.pytorch.ddpg import DDPG
10	10	- from utils.env import launch_env
11	11	- from utils.wrappers import NormalizeWrapper, ImgWrapper, \ + from learning.reinforcement.pytorch.ddpg import DDPG + from learning.utils.env import launch_env + from learning.utils.wrappers import NormalizeWrapper, ImgWrapper, \ 12     DtRewardWrapper, ActionWrapper, ResizeWrapper
13	13	
14	14	

▼ BIN +19.6 MB learning/reinforcement/pytorch/models/planet.pth/model.pth ↗

Binary file not shown.

▼ 193 learning/reinforcement/pytorch/test\_algo.py ↗

...	...	@@ -0,0 +1,193 @@
	1	+ ##### + ##### FOR GLEN ##### + ##### 4   + # This is just a script I used to help figure out how to record videos in 5   + # Duckietown, and also how to use a trained agent to pick actions in Duckietown. 6   + ##### 7   + import ast 8   + import argparse 9   + import logging 10   + 11   + import os 12   + from tkinter import W 13   + import numpy as np 14   + 15   + from gym.wrappers import Monitor 16   + from gym.wrappers.monitoring.video_recorder import VideoRecorder 17   + 18   + # Duckietown Specific 19   + from learning.reinforcement.pytorch.ddpg import DDPG 20   + from learning.utils.env import launch_env 21   + from learning.utils.wrappers import NormalizeWrapper, ImgWrapper, \ 22       DtRewardWrapper, ActionWrapper, ResizeWrapper 23   + 24   + import wandb 25   + import torch 26   + 27   + import errno 28   + import os 29   + from datetime import datetime 30   + 31   + import hydra 32   + import omegaconf

```

33 + from learning.reinforcement.pytorch.utils import planet_config_dict
34 +
35 + cfg = omegaconf.OmegaConf.create(planet_config_dict)
36 +
37 + # ADD THE MBRL LIB TO PATH SO YOU CAN IMPORT IT
38 + import sys
39 + sys.path.append('/home/kuwajerw/repos/duckietown-mbrl-lib')
40 +
41 +
42 + import mbrl.algorithms.planet as planet
43 + import mbrl.util.env
44 + import mbrl.constants
45 + from mbrl.env.termination_fns import no_termination
46 + from mbrl.models import ModelEnv, ModelTrainer
47 + from mbrl.planning import RandomAgent, create_trajectory_optim_agent_for_model
48 +
49 + # folder in which to save logs
50 + # will make timestamped folder in here
51 + base_logdir = 'saved_runs'
52 +
53 + # timestamp stuff
54 + base_day_str = datetime.now().strftime('%Y-%m-%d')
55 + base_time_str = datetime.now().strftime('%H-%M-%S')
56 + curr_timestamp_str = base_day_str + '_' + base_time_str
57 +
58 + # made this helper function to manually save the logs
59 + # in addition to the W&B logs just in case
60 + def make_save_folder(label=''):
61 +
62 +     # make timestamp logdir path
63 +     base_path = os.getcwd() + f'/{base_logdir}/'
64 +     mydir = os.path.join(
65 +         base_path,
66 +         base_day_str, base_time_str + "_" + label)
67 +
68 +     # actually create the folder
69 +     try:
70 +         os.makedirs(mydir)
71 +     except OSError as e:
72 +         if e.errno != errno.EEXIST:
73 +             raise # This was not a "directory exist" error..
74 +
75 +     return mydir + "/"
76 +
77 +
78 + def save_list(items, label, dest):
79 +     with open(os.path.join(dest, label), 'w') as f:
80 +         f.writelines(items)
81 +
82 +
83 +
84 + # curr_timestamp_str = datetime.now().strftime('%Y-%m-%d_%H-%M-%S')
85 +
86 + video_name = f"{curr_timestamp_str}.mp4"
87 +
88 +
89 +
90 + def _enjoy():
91 +
92 +     # save run logs in folder that matches the wandb run name
93 +     logdir = make_save_folder(label=wandb.run.name)
94 +
95 +     # Launch the env with our helper function
96 +     env = launch_env()
97 +     print("Initialized environment")
98 +
99 +     # Wrappers
100

```

```

101 +     env = ResizeWrapper(env)
102 +     env = NormalizeWrapper(env)
103 +     env = ImgWrapper(env) # to make the images from 160x120x3 into 3x160x120
104 +     env = ActionWrapper(env)
105 +     env = DtRewardWrapper(env)
106 +     print("Initialized Wrappers")
107 +
108 +     # record video
109 +     video_rec = VideoRecorder(env, logdir + video_name)
110 +
111 +     state_dim = env.observation_space.shape
112 +     action_dim = env.action_space.shape[0]
113 +     max_action = float(env.action_space.high[0])
114 +
115 +     # PLANET STUFF #####
116 +
117 +     obs_shape = env.observation_space.shape
118 +     act_shape = env.action_space.shape
119 +     # Create PlaNet model
120 +     rng = torch.Generator(device=cfg.device)
121 +     rng.manual_seed(cfg.seed)
122 +     cfg.dynamics_model.action_size = env.action_space.shape[0]
123 +
124 +     planet = hydra.utils.instantiate(cfg.dynamics_model)
125 +     assert isinstance(planet, mbrl.models.PlaNetModel)
126 +     planet.load('learning/reinforcement/pytorch/models/planet.pth')
127 +     model_env = ModelEnv(env, planet, no_termination, generator=rng)
128 +     trainer = ModelTrainer(planet, optim_lr=1e-3, optim_eps=1e-4)
129 +
130 +     agent = create_trajectory_optim_agent_for_model(model_env, cfg.algorithm.agent)
131 +     # PLANET STUFF #####
132 +
133 +     # Initialize policy
134 +     # policy = DDPG(state_dim, action_dim, max_action, net_type="cnn")
135 +     # policy.load(filename='ddpg', directory='learning/reinforcement/pytorch/models/')
136 +
137 +     obs = env.reset()
138 +     done = False
139 +     reward_list = []
140 +     action_list = []
141 +
142 +     # PLANET STUFF #####
143 +     agent.reset()
144 +     planet.reset_posterior()
145 +     action = None
146 +     # PLANET STUFF #####
147 +
148 +     # while True:
149 +     # while not done:
150 +     for step in range(100):
151 +
152 +         # PLANET STUFF #####
153 +         planet.update_posterior(obs, action=action, rng=rng)
154 +         action_noise = 0
155 +         action = agent.act(obs) + action_noise
156 +         action = np.clip(action, -1.0, 1.0) # to account for the noise
157 +         # PLANET STUFF #####
158 +
159 +         # action = policy.predict(np.array(obs))
160 +         # action = env.action_space.sample()
161 +         # Perform action
162 +         obs, reward, done, _ = env.step(action)
163 +
164 +         # log stuff with wandb
165 +         wandb.log({'step' : step})
166 +         wandb.log({'action0': action[0], 'action1' : action[1], 'global_step': step})
167

```

```

168 +         wandb.log({'reward' : reward, 'global_step': step})
169 +
170 +     # env.render() # optional
171 +     video_rec.capture_frame() # record video
172 +
173 +     # maintain manual log just in case
174 +     reward_list.append(str(reward) + '\n')
175 +     action_list.append(str(action) + '\n')
176 +
177 +
178 +     # cleanup here
179 +     # done = False
180 +     # obs = env.reset()
181 +     video_rec.close()
182 +     env.close()
183 +
184 +     # save manual log
185 +     save_list(reward_list, f"rewards_{curr_timestamp_str}.txt", logdir)
186 +     save_list(action_list, f"actions_{curr_timestamp_str}.txt", logdir)
187 +
188 +
189 +
190 + if __name__ == '__main__':
191 +     wandb.init(project="MBRL_Duckyt", entity="mbrl_ducky", monitor_gym=True)
192 +     _enjoy()
193 +

```

▼ ⇤ 10 learning/reinforcement/pytorch/train\_reinforcement.py □

6	6	import numpy as np
7	7	
8	8	# Duckietown Specific
9	9	- from reinforcement.pytorch.ddpg import DDPG
10	10	- from reinforcement.pytorch.utils import seed, evaluate_policy, ReplayBuffer
11	11	- from utils.env import launch_env
12	12	- from utils.wrappers import NormalizeWrapper, ImgWrapper, \ + from learning.reinforcement.pytorch.ddpg import DDPG + from learning.reinforcement.pytorch.utils import seed, evaluate_policy, ReplayBuffer + from learning.utils.env import launch_env + from learning.utils.wrappers import NormalizeWrapper, ImgWrapper, \ + DtRewardWrapper, ActionWrapper, ResizeWrapper
13	13	
14	14	
15	15	logger = logging.getLogger(__name__)
138	138	parser.add_argument("--policy_freq", default=2, type=int) # Frequency of delayed policy updates
139	139	parser.add_argument("--env_timesteps", default=500, type=int) # Frequency of delayed policy updates
140	140	parser.add_argument("--replay_buffer_max_size", default=10000, type=int) # Maximum number of steps to keep
141	141	- parser.add_argument('--model-dir', type=str, default='reinforcement/pytorch/models/')
141	141	+ parser.add_argument('--model-dir', type=str, default='learning/reinforcement/pytorch/models/')
142	142	
143	143	_train(parser.parse_args())

▼ ⇤ 4 manual\_control.py □

38	38	domain_rand = args.domain_rand,
39	39	frame_skip = args.frame_skip,
40	40	distortion = args.distortion,
41	41	+ set_start_pos=np.array([3.68356218, 0., 1.50287902]),
42	42	+ set_start_angle=np.array([-1.564463624086557])
41	43	)
42	44	else:
43	45	env = gym.make(args.env_name)
98	100	
99	101	obs, reward, done, info = env.step(action)
100	102	print('step_count = %s, reward=%3f' % (env.unwrapped.step_count, reward))
103	103	+ print(f"curr_pos is {env.cur_pos} and cur_angle is {env.cur_angle}")
104	104	+

101	105
102	106
103	107

```
if key_handler[key.RETURN]:  
    from PIL import Image
```