🗄 Coac / **never-give-up**   Public

`<>` **Code**   ⊙ Issues  1   ⑂ Pull requests   ▷ Actions   ▦ Projects   ⛨ Security   ···

⑂ main ▾                                                         ···

**never-give-up** / **train.py** / `<>`Jump to ▾

🖼 **Coac** Add working config                          ⟲ **History**

👥 **1** contributor

161 lines (124 sloc)  │  5.11 KB                                    ···

```python
1    import numpy as np
2    import torch
3    import torch.optim as optim
4    import wandb
5    from gym import Wrapper
6    from gym_maze.envs.maze_env import MazeEnvSample5x5
7
8    from config import config
9    from embedding_model import EmbeddingModel, compute_intrinsic_reward
10   from memory import Memory, LocalBuffer
11   from model import R2D2
12
13
14   def get_action(state, target_net, epsilon, env, hidden):
15       action, hidden = target_net.get_action(state, hidden)
16
17       if np.random.rand() <= epsilon:
18           return env.action_space.sample(), hidden
19       else:
20           return action, hidden
21
22
23   def update_target_model(online_net, target_net):
24       target_net.load_state_dict(online_net.state_dict())
25
26
27   class Maze(Wrapper):
28       def step(self, action: int):
29           obs, rew, done, info = super().step(["N", "E", "S", "W"][action])
30           self.set.add((obs[0], obs[1]))
31           if rew > 0:
32               rew = 10
```

*Handwritten annotations:*
- (near line 14) $n_t$ in paper
- (near line 14) from R2D2 code, produced by LSTM → handles POMDP
- (near lines 17–20) → ε-greedy
- (near line 24) ← straight copy of online network weights
- (near line 30) Move the agent

```python
33            return obs / 10, rew, done, info
34
35        def reset(self):
36            self.set = set()
37            return super().reset()
38
39
40   def main():
41        env = Maze(MazeEnvSample5x5())
42
43        torch.manual_seed(config.random_seed)
44        env.seed(config.random_seed)
45        np.random.seed(config.random_seed)
46        env.action_space.seed(config.random_seed)
47
48        wandb.init(project="ngu-maze", config=config.__dict__)
49
50        num_inputs = env.observation_space.shape[0]
51        num_actions = env.action_space.n
52        print("state size:", num_inputs)
53        print("action size:", num_actions)
54
55        online_net = R2D2(num_inputs, num_actions)
56        target_net = R2D2(num_inputs, num_actions)
57        update_target_model(online_net, target_net)
58        embedding_model = EmbeddingModel(obs_size=num_inputs, num_outputs=num_actions)
59        embedding_loss = 0
60
61        optimizer = optim.Adam(online_net.parameters(), lr=config.lr)
62
63        online_net.to(config.device)
64        target_net.to(config.device)
65        online_net.train()
66        target_net.train()
67        memory = Memory(config.replay_memory_capacity)
68        epsilon = 1.0
69        steps = 0
70        loss = 0
71        local_buffer = LocalBuffer()
72        sum_reward = 0
73        sum_augmented_reward = 0
74        sum_obs_set = 0
75
76        for episode in range(30000):
77            done = False
78            state = env.reset()
79            state = torch.Tensor(state).to(config.device)
80
81            hidden = (
82                torch.Tensor().new_zeros(1, 1, config.hidden_size),
83                torch.Tensor().new_zeros(1, 1, config.hidden_size),
```

Handwritten annotations:
- (line 33) termination flag → done
- (line 34) state → obs; reward → rew
- (lines 43–46) Note for reproducibility. Set random seeds within every class
- (line 54) see model.py
- (line 57) Same weights at the start
- (line 59) embedding-model.py → see left side of fig 1 in paper
- (line 76) number of episodes

```python
84                   )
85
86               episodic_memory = [embedding_model.embedding(state)]
87
88               episode_steps = 0
89               horizon = 100
90               while not done:
91                   steps += 1
92                   episode_steps += 1
93
94                   action, new_hidden = get_action(state, target_net, epsilon, env, hidden)
95
96                   next_state, env_reward, done, _ = env.step(action)
97                   next_state = torch.Tensor(next_state)
98
99                   augmented_reward = env_reward
100                  if config.enable_ngu:
101                      next_state_emb = embedding_model.embedding(next_state)
102                      intrinsic_reward = compute_intrinsic_reward(episodic_memory, next_st
103                      episodic_memory.append(next_state_emb)
104                      beta = 0.0001
105                      augmented_reward = env_reward + beta * intrinsic_reward
106
107                  mask = 0 if done else 1
108
109                  local_buffer.push(state, next_state, action, augmented_reward, mask, hid
110                  hidden = new_hidden
111                  if len(local_buffer.memory) == config.local_mini_batch:
112                      batch, lengths = local_buffer.sample()
113                      td_error = R2D2.get_td_error(online_net, target_net, batch, lengths)
114                      memory.push(td_error, batch, lengths)
115
116                  sum_reward += env_reward
117                  state = next_state
118                  sum_augmented_reward += augmented_reward
119
120                  if steps > config.initial_exploration and len(memory) > config.batch_siz
121                      epsilon -= config.epsilon_decay
122                      epsilon = max(epsilon, 0.4)
123
124                      batch, indexes, lengths = memory.sample(config.batch_size)
125                      loss, td_error = R2D2.train_model(online_net, target_net, optimizer,
126
127                      if config.enable_ngu:
128                          embedding_loss = embedding_model.train_model(batch)
129
130                      memory.update_priority(indexes, td_error, lengths)
131
132                      if steps % config.update_target == 0:
133                          update_target_model(online_net, target_net)
134
```

*(handwritten annotations)*

replace

running in NGU mode

→ see embedding-model.py

↳ RL training

↳ training embedding

slow update after fixed number of updates of online model

```python
135                  if episode_steps >= horizon or done:
136                      sum_obs_set += len(env.set)
137                      break
138
139              if episode > 0 and episode % config.log_interval == 0:
140                  mean_reward = sum_reward / config.log_interval
141                  mean_augmented_reward = sum_augmented_reward / config.log_interval
142                  metrics = {
143                      "episode": episode,
144                      "mean_reward": mean_reward,
145                      "epsilon": epsilon,
146                      "embedding_loss": embedding_loss,
147                      "loss": loss,
148                      "mean_augmented_reward": mean_augmented_reward,
149                      "steps": steps,
150                      "sum_obs_set": sum_obs_set / config.log_interval,
151                  }
152                  print(metrics)
153                  wandb.log(metrics)
154
155                  sum_reward = 0
156                  sum_augmented_reward = 0
157                  sum_obs_set = 0
158
159
160     if __name__ == "__main__":
161         main()
```

*Useful for tracking if you haven't tried it* (handwritten annotation)