

# DQN-Pong Review

## 雅达利游戏智能体设计报告

康振飞，张润泽

云南大学-新技术专题 [1]

2024 年 12 月 20 日



## ① 课题背景

## ② 研究现状

## ③ 研究内容

## ④ 训练结果

## ① 课题背景

## ② 研究现状

## ③ 研究内容

## ④ 训练结果

# 强化学习基础研究

- 验证算法性能:

# 强化学习基础研究

- 验证算法性能：
  - Atari Pong 是一个经典的强化学习基准任务，训练智能体玩 Pong 可以验证 DQN 算法的性能和鲁棒性。通过在 Pong 上的实验，可以测试和改进 DQN 的超参数、网络结构、经验回放机制等。

# 强化学习基础研究

- 验证算法性能：
  - Atari Pong 是一个经典的强化学习基准任务，训练智能体玩 Pong 可以验证 DQN 算法的性能和鲁棒性。通过在 Pong 上的实验，可以测试和改进 DQN 的超参数、网络结构、经验回放机制等。
- 算法优化：

# 强化学习基础研究

- 验证算法性能：
  - Atari Pong 是一个经典的强化学习基准任务，训练智能体玩 Pong 可以验证 DQN 算法的性能和鲁棒性。通过在 Pong 上的实验，可以测试和改进 DQN 的超参数、网络结构、经验回放机制等。
- 算法优化：
  - 通过实验可以发现 DQN 的瓶颈（如样本效率低、训练不稳定等），从而为后续算法的改进提供方向。

# 游戏智能体应用

- 游戏 AI 研究：



# 游戏智能体应用

- 游戏 AI 研究：
  - 通过训练 Pong 智能体，可以研究游戏 AI 的行为模式、策略选择和学习能力。

# 游戏智能体应用

- 游戏 AI 研究：
  - 通过训练 Pong 智能体，可以研究游戏 AI 的行为模式、策略选择和学习能力。
- 辅助理解架构知识：

# 游戏智能体应用

- 游戏 AI 研究：
  - 通过训练 Pong 智能体，可以研究游戏 AI 的行为模式、策略选择和学习能力。
- 辅助理解架构知识：
  - 通过调整超参和修改局部代码，可以进一步理解强化学习的概念（如状态、动作、奖励、Q 值等）。

## ① 课题背景

## ② 研究现状

## ③ 研究内容

## ④ 训练结果

# 基础网络架构<sup>1</sup>

- 输入层：Input 是一个  $84*84*4$  的图像。使用了最近的 4 帧画面，每帧  $84*84$  像素。
- 卷积层：
  - 1st: 32 个  $8*8$  滤波器，步幅 4。从输入图像中提取基础特征，如边缘和简单形状。
  - 2nd: 64 个  $4*4$  滤波器，步幅 2。这一层提取更复杂的特征，步幅减小以获得更详细的特征图。
  - 3rd: 64 个  $3*3$  滤波器，步幅 1。这一层进一步提取精细特征，步幅更小以保持特征的细节。
- 全链接层：
  - 1st: `nn.Linear(64*7*7, 512)`。输入特征数：3136，输出特征数：512（神经元数）。
  - ReLU：负值置 0，正值不变。
  - 2nd: `nn.Linear(512, action-space.n)`。输入特征数：512（接上层），输出至动作空间。

---

<sup>1</sup>模型参考 doi:10.1038/nature14236

## ① 课题背景

## ② 研究现状

## ③ 研究内容

网络剖析

训练及超参设置

## ④ 训练结果

## ① 课题背景

## ② 研究现状

## ③ 研究内容

网络剖析

训练及超参设置

## ④ 训练结果

# DQN 设计<sup>2</sup>

```

01. class DQN(nn.Module):
02.     def __init__(self,
03.         observation_space: spaces.Box,
04.         action_space: spaces.Discrete):
05.         super().__init__()
06.         assert type(
07.             observation_space) == spaces.Box, 'observation_space must be of type Box'
08.         assert len(
09.             observation_space.shape) == 3, 'observation space must have the form channels x width x height'
10.         assert type(
11.             action_space) == spaces.Discrete, 'action_space must be of type Discrete'
12.
13.         self.conv = nn.Sequential(
14.             nn.Conv2d(in_channels=observation_space.shape[0], out_channels=32, kernel_size=8, stride=4),
15.             nn.ReLU(),
16.             nn.Conv2d(in_channels=32, out_channels=64, kernel_size=4, stride=2),
17.             nn.ReLU(),
18.             nn.Conv2d(in_channels=64, out_channels=64, kernel_size=3, stride=1),
19.             nn.ReLU()
20.         )
21.
22.         self.fc = nn.Sequential(
23.             nn.Linear(in_features=64*7*7, out_features=512),
24.             nn.ReLU(),
25.             nn.Linear(in_features=512, out_features=action_space.n)
26.         )
27.
28.     def forward(self, x):
29.         conv_out = self.conv(x).view(x.size()[0],-1)
30.         return self.fc(conv_out)
    
```

图 1: DQN Code

<sup>2</sup>完整代码: [github.com/godsboy404/DQN-Pong](https://github.com/godsboy404/DQN-Pong)



图 2: LifeCycle & Rewards

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻



## ① 课题背景

## ② 研究现状

## ③ 研究内容

网络剖析

训练及超参设置

## ④ 训练结果

## 超参

```
01. hyper_params = {
02.     "env": "PongNoFrameskip-v4", # name of the game
03.     "replay-buffer-size": int(5e3), # replay buffer size
04.     "learning-rate": 1e-4, # learning rate for Adam optimizer
05.     "discount-factor": 0.99, # discount factor
06.     "num-steps": int(1e6), # total number of steps to run the environment for
07.     "batch-size": 32, # number of transitions to optimize at the same time
08.     "learning-starts": 10000, # number of steps before learning starts
09.     "learning-freq": 1, # number of iterations between every optimization step
10.     "use-double-dqn": True, # use double deep Q-learning
11.     "target-update-freq": 1000, # number of iterations between every target network update
12.     "eps-start": eps_start, # e-greedy start threshold
13.     "eps-end": 0.01, # e-greedy end threshold
14.     "eps-fraction": 0.1, # fraction of num-steps
15.     "print-freq": 10
16. }
```

图 4: params

图 5: train

```

02. state = env.reset()
03. eps_timesteps = hyper_params["eps-fraction"] * \
04.     float(hyper_params["num-steps"])
05. episode_rewards = [0,0]
06.
07. for t in range(hyper_params["num-steps"]):
08.     fraction = min(1.0, float(t) / eps_timesteps)
09.     eps_threshold = hyper_params["eps-start"] + fraction * \
10.         (hyper_params["eps-end"] - hyper_params["eps-start"])
11.     sample = randoms.random()
12.
13.     if(sample > eps_threshold):
14.         # Exploit
15.         action = agent.act(state)
16.     else:
17.         # Explore
18.         action = env.action_space.sample()
19.
20.     next_state, reward, done, info = env.step(action)
21.     agent.memory.add(state, action, reward, next_state, done)
22.     state = next_state
23.
24.     episode_rewards[-1] += reward
25.     if done:
26.         state = env.reset()
27.         episode_rewards.append(0,0)
28.
29.     if t > hyper_params["learning-starts"] and t % hyper_params["learning-freq"] == 0:
30.         agent.optimise_td_loss()
31.
32.     if t > hyper_params["learning-starts"] and t % hyper_params["target-update-freq"] == 0:
33.         agent.update_target_network()
34.
35. num_episodes = len(episode_rewards)
36.
37. if done and hyper_params["print-freq"] is not None and len(episode_rewards) % hyper_params["print-freq"] == 0:
38.     mean_100ep_reward = round(np.mean(episode_rewards[-101:-1]), 1)
39.     print("\tEpisode: {}".format(t))
40.     print("\tEpisodes: {}".format(num_episodes))
41.     print("\tMean 100 episode reward: {}".format(mean_100ep_reward))
42.     print("\t% time spent exploring: {}".format(int(100 * eps_threshold)))
43.     print(".....")
44.     torch.save(agent.policy_network.state_dict(), f'checkpoint_{t}.pth')
45.     np.savetxt('results/rewards_per_episode.csv', episode_rewards,
46.               delimiter=',', fmt='%1.3f')

```

# 优化误差

```
01. def optimise_td_loss(self):
02.     """
03.     Optimise the TD-error over a single minibatch of transitions
04.     :return: the loss
05.     """
06.     device = self.device
07.
08.     states, actions, rewards, next_states, dones = self.memory.sample(self.batch_size)
09.     states = np.array(states) / 255.0
10.     next_states = np.array(next_states) / 255.0
11.     states = torch.from_numpy(states).float().to(device)
12.     actions = torch.from_numpy(actions).long().to(device)
13.     rewards = torch.from_numpy(rewards).float().to(device)
14.     next_states = torch.from_numpy(next_states).float().to(device)
15.     dones = torch.from_numpy(dones).float().to(device)
16.
17.     with torch.no_grad():
18.         if self.use_double_dqn:
19.             _, max_next_action = self.policy_network(next_states).max(1)
20.             max_next_q_values = self.target_network(next_states).gather(1, max_next_action.unsqueeze(1)).squeeze()
21.         else:
22.             next_q_values = self.target_network(next_states)
23.             max_next_q_values, _ = next_q_values.max(1)
24.             target_q_values = rewards + (1 - dones) * self.gamma * max_next_q_values
25.
26.     input_q_values = self.policy_network(states)
27.     input_q_values = input_q_values.gather(1, actions.unsqueeze(1)).squeeze()
28.
29.     loss = F.smooth_l1_loss(input_q_values, target_q_values)
30.
31.     self.optimiser.zero_grad()
32.     loss.backward()
33.     self.optimiser.step()
34.     del states
35.     del next_states
36.     return loss.item()
```

图 6: loss

# 采取动作 & 更新陪练网络

```
01. def act(self, state: np.ndarray):
02.     """
03.     Select an action greedily from the Q-network given the state
04.     :param state: the current state
05.     :return: the action to take
06.     """
07.     device = self.device
08.     state = np.array(state) / 255.0
09.     state = torch.from_numpy(state).float().unsqueeze(0).to(device)
10.     with torch.no_grad():
11.         q_values = self.policy_network(state)
12.         _, action = q_values.max(1)
13.         return action.item()
14.
15. def update_target_network(self):
16.     """
17.     Update the target Q-network by copying the weights from the current Q-network
18.     """
19.     self.target_network.load_state_dict(self.policy_network.state_dict())
```

图 7: action & network-updating

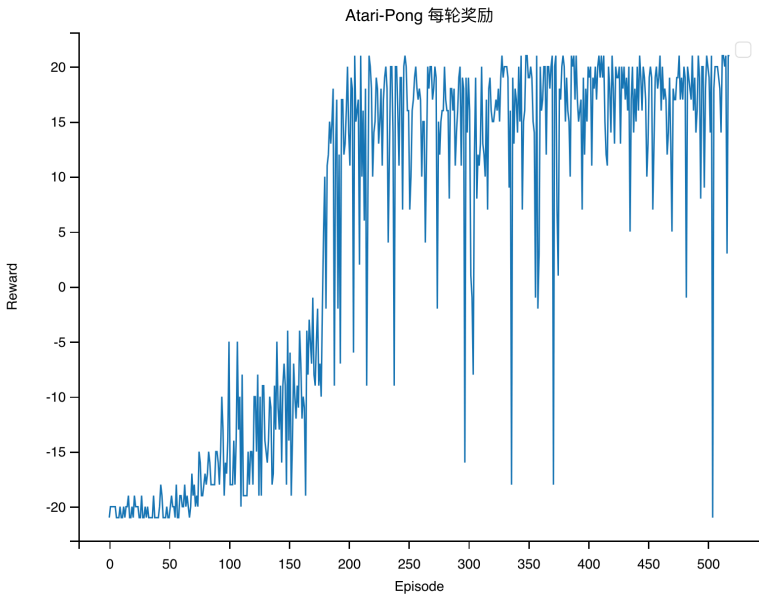
## ① 课题背景

## ② 研究现状

## ③ 研究内容

## ④ 训练结果





*Thanks for Watching!*