



Федеральное государственное бюджетное образовательное учреждение
высшего профессионального образования
«Московский государственный технический университет имени Н.Э.
Баумана»
(МГТУ им. Н.Э. Баумана)

Лабораторная работа №6
по курсу «Методы машинного обучения»

Выполнил
студент группы ИУ5-22М
XXXX

Москва, 2023

1. Задание

- На основе рассмотренных на лекции примеров реализуйте алгоритм DQN.
- В качестве среды можно использовать классические среды (в этом случае используется полносвязная архитектура нейронной сети).
- В качестве среды можно использовать игры Atari (в этом случае используется сверточная архитектура нейронной сети).
- В случае реализации среды на основе сверточной архитектуры нейронной сети +1 балл за экзамен.

2. Текст программы

```
1 #!/usr/bin/env python
2
3 import gymnasium as gym
4 import math
5 import random
6 import matplotlib.pyplot as plt
7 from collections import namedtuple, deque
8 import torch
9 import torch.nn as nn
10 import torch.optim as optim
11 import torch.nn.functional as F
12
13 # Название среды
14 CONST_ENV_NAME = 'Acrobot-v1'
15
16 # Использование GPU
17 CONST_DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
18
19 # Элемент ReplayMemory в форме именованного кортежа
20 Transition = namedtuple('Transition', ('state', 'action', 'next_state', 'reward'))
21
22 # Реализация техники Replay Memory
23 class ReplayMemory(object):
24     def __init__(self, capacity):
25         self.memory = deque([], maxlen=capacity)
26
27     def push(self, *args):
28         '''
29         Сохранение данных в ReplayMemory
30         '''
31
32         self.memory.append(Transition(*args))
33
34     def sample(self, batch_size):
35         '''
36         Выборка случайных элементов размера batch_size
37         '''
38
39         return random.sample(self.memory, batch_size)
40
41     def __len__(self):
42         return len(self.memory)
43
44 class DQN_Model(nn.Module):
```

```

45 def __init__(self, n_observations, n_actions):
46     '''
47     Инициализация топологии нейронной сети
48     '''
49
50     super(DQN_Model, self).__init__()
51     self.layer1 = nn.Linear(n_observations, 128)
52     self.layer2 = nn.Linear(128, 64)
53     self.layer3 = nn.Linear(64, n_actions)
54
55 def forward(self, x):
56     '''
57     Прямой проход
58     Вызывается для одного элемента, чтобы определить следующее действие
59     Или для batch во время процедуры оптимизации
60     '''
61
62     x = F.relu(self.layer1(x))
63     x = F.relu(self.layer2(x))
64     return self.layer3(x)
65
66 class DQN_Agent:
67     def __init__(
68         self,
69         env,
70         BATCH_SIZE = 128,
71         GAMMA = 0.99,
72         EPS_START = 0.1,
73         EPS_END = 0.5,
74         EPS_DECAY = 1000,
75         TAU = 0.005,
76         LR = 0.0001,
77     ):
78         # Среда
79         self.env = env
80         # Размерности Q-модели
81         self.n_actions = env.action_space.n
82         state, _ = self.env.reset()
83         self.n_observations = len(state)
84         # Коэффициенты
85         self.BATCH_SIZE = BATCH_SIZE
86         self.GAMMA = GAMMA
87         self.EPS_START = EPS_START
88         self.EPS_END = EPS_END
89         self.EPS_DECAY = EPS_DECAY
90         self.TAU = TAU
91         self.LR = LR
92
93         # Модели
94         # Основная модель
95         self.policy_net = DQN_Model(self.n_observations, self.n_actions).to(CONST_DEVICE)
96
97         # Вспомогательная модель, используется для стабилизации алгоритма
98         # Обновление контролируется гиперпараметром TAU
99         # Используется подход Double DQN
100        self.target_net = DQN_Model(self.n_observations, self.n_actions).to(CONST_DEVICE)
101        self.target_net.load_state_dict(self.policy_net.state_dict())
102
103        # Оптимизатор
104        self.optimizer = optim.AdamW(self.policy_net.parameters(), lr=self.LR, amsgrad=True)
105

```

```

106 # Replay Memory
107 self.memory = ReplayMemory(10000)
108
109 # Количество шагов
110 self.steps_done = 0
111
112 # Длительность эпизодов
113 self.episode_durations = []
114
115 def select_action(self, state):
116     '''
117     Выбор действия
118     '''
119
120     sample = random.random()
121     eps = self.EPS_END + (self.EPS_START - self.EPS_END) * math.exp(-1. * self.steps_done /
122 self.EPS_DECAY)
123     self.steps_done += 1
124     if sample > eps:
125         with torch.no_grad():
126             # Если вероятность больше eps
127             # то выбирается действие, соответствующее максимальному Q-значению
128             # t.max(1) возвращает максимальное значение колонки для каждой строки
129             # [1] возвращает индекс максимального элемента
130             return self.policy_net(state).max(1)[1].view(1, 1)
131     else:
132         # Если вероятность меньше eps
133         # то выбирается случайное действие
134         return torch.tensor([self.env.action_space.sample()], device=CONST_DEVICE,
135 dtype=torch.long)
136
137 def plot_durations(self, show_result=False):
138     plt.figure(1)
139     durations_t = torch.tensor(self.episode_durations, dtype=torch.float)
140     if show_result:
141         plt.title('Результат')
142     else:
143         plt.clf()
144         plt.title('Обучение')
145         plt.xlabel('Эпизод')
146         plt.ylabel('Количество шагов в эпизоде')
147         plt.plot(durations_t.numpy())
148         plt.pause(0.001) # пауза
149
150 def optimize_model(self):
151     '''
152     Оптимизация модели
153     '''
154
155     if len(self.memory) < self.BATCH_SIZE:
156         return
157
158     transitions = self.memory.sample(self.BATCH_SIZE)
159     # Транспонирование batch'a
160     # Конвертация batch-массива из Transition
161     # в Transition batch-массивов.
162     batch = Transition(*zip(*transitions))
163
164     # Вычисление маски нефинальных состояний и конкатенация элементов batch'a
165     non_final_mask = torch.tensor(tuple(map(lambda s: s is not None, batch.next_state)),
166 device=CONST_DEVICE, dtype=torch.bool)

```

```

164 non_final_next_states = torch.cat([s for s in batch.next_state if s is not None])
165 state_batch = torch.cat(batch.state)
166 action_batch = torch.cat(batch.action)
167 reward_batch = torch.cat(batch.reward)
168
169 # Вычисление  $Q(s_t, a)$ 
170 state_action_values = self.policy_net(state_batch).gather(1, action_batch)
171
172 # Вычисление  $V(s_{t+1})$  для всех следующих состояний
173 next_state_values = torch.zeros(self.BATCH_SIZE, device=CONST_DEVICE)
174
175 with torch.no_grad():
176     next_state_values[non_final_mask] = self.target_net(non_final_next_states).max(1)[0]
177
178 # Вычисление ожидаемых значений  $Q$ 
179 expected_state_action_values = (next_state_values * self.GAMMA) + reward_batch
180
181 # Вычисление Huber loss
182 criterion = nn.SmoothL1Loss()
183 loss = criterion(state_action_values, expected_state_action_values.unsqueeze(1))
184
185 # Оптимизация модели
186 self.optimizer.zero_grad()
187 loss.backward()
188
189 # gradient clipping
190 torch.nn.utils.clip_grad_value_(self.policy_net.parameters(), 100)
191 self.optimizer.step()
192
193 def play_agent(self):
194     '''
195     Проигрывание сессии для обученного агента
196     '''
197
198     env2 = gym.make(CONST_ENV_NAME, render_mode='human')
199     state = env2.reset()[0]
200     state = torch.tensor(state, dtype=torch.float32, device=CONST_DEVICE).unsqueeze(0)
201     res = []
202
203     terminated = False
204     truncated = False
205
206     while not terminated and not truncated:
207         action = self.select_action(state)
208         action = action.item()
209         observation, reward, terminated, truncated, _ = env2.step(action)
210         env2.render()
211         res.append((action, reward))
212
213         state = torch.tensor(observation, dtype=torch.float32, device=CONST_DEVICE).unsqueeze(0)
214
215     print('done!')
216     print('Данные об эпизоде: ', res)
217
218 def train(self):
219     '''
220     Обучение агента
221     '''
222
223     if torch.cuda.is_available():
224         num_episodes = 600

```

```

225     else:
226         num_episodes = 50
227
228     for i_episode in range(num_episodes):
229         # Инициализация среды
230         state, info = self.env.reset()
231         state = torch.tensor(state, dtype=torch.float32, device=CONST_DEVICE).unsqueeze(0)
232
233         terminated = False
234         truncated = False
235
236         iters = 0
237         while not terminated and not truncated:
238             action = self.select_action(state)
239             observation, reward, terminated, truncated, _ = self.env.step(action.item())
240             reward = torch.tensor([reward], device=CONST_DEVICE)
241
242             if terminated:
243                 next_state = None
244             else:
245                 next_state = torch.tensor(observation, dtype=torch.float32,
246                                           device=CONST_DEVICE).unsqueeze(0)
247
248             # Сохранение данных в Replay Memory
249             self.memory.push(state, action, next_state, reward)
250
251             # Переход к следующему состоянию
252             state = next_state
253
254             # Выполнение одного шага оптимизации модели
255             self.optimize_model()
256
257             # Обновление весов target-сети
258             #  $\theta' \leftarrow \tau \theta + (1 - \tau) \theta'$ 
259             target_net_state_dict = self.target_net.state_dict()
260             policy_net_state_dict = self.policy_net.state_dict()
261
262             for key in policy_net_state_dict:
263                 target_net_state_dict[key] = policy_net_state_dict[key] * self.TAU +
264                 target_net_state_dict[key] * (1 - self.TAU)
265
266             self.target_net.load_state_dict(target_net_state_dict)
267             iters += 1
268
269             self.episode_durations.append(iters)
270             self.plot_durations()
271
272 def main():
273     env = gym.make(CONST_ENV_NAME)
274     agent = DQN_Agent(env)
275     agent.train()
276     agent.play_agent()
277
278 if __name__ == '__main__':
279     main()

```

3. Экранные формы с примерами выполнения программы

