

A true man can play a palo one hundred time

发表于 2021-01-21 分类于 [Challenge](#) , [2017](#) , [HCTF](#) , [Web](#)
Challenge | 2017 | HCTF | Web | A true man can play a palo one hundred time

[点击此处](#)获得更好的阅读体验

WriteUp 来源

<https://xz.aliyun.com/t/1589>

题目考点

解题思路

题目说明

```
1 Question
2 Now you have a balance palo.
3 You can move it left or right.
4 Just play hundred time on it.
5
6 Description
7 Get request receive two params
8   1. move, 0 or 1
9   2. id, just your token
10 Observation
11
12   1. pole position x
13   2. a value depend on x
14   3. pole deviate from center f.
15   4. a value depend on f.
16 Why you failed
17 f. or x > a certain value
```

总而言之就是个玩棒子的游戏(雾。之所以出现在最后一道请去问关卡规则的设计者。因为ctf本来不应该出现这种问题，所以我有意把这题设计得简单了一点，但是，ctf真是不讲道理，也导致这道题被少量非预期。

其实就是一个非常非常简单的强化学习的问题，甚至不需要强化学习去解。

DQN网络结构定义

```
1 import numpy as np
2 import tensorflow as tf
3 import requests
4 import math
5
6 class DeepQNetwork:
7     def __init__(self,
8                  n_actions,
9                  n_features,
10                 learning_rate=0.01,
11                 reward_decay=0.9,
12                 e_greedy=0.9,
13                 replace_target_iter=300,
14                 memory_size=500,
15                 batch_size=32,
16                 e_greedy_increment=None,
17                 output_graph=False,
18                 ):
19         self.n_actions = n_actions
20         self.n_features = n_features
21         self.lr = learning_rate
22         self.gamma = reward_decay
23         self.epsilon_max = e_greedy
24         self.replace_target_iter = replace_target_iter
25         self.memory_size = memory_size
26         self.batch_size = batch_size
27         self.epsilon_increment = e_greedy_increment
28         self.epsilon = 0 if e_greedy_increment is not None else self.epsilon_max
29
30         # total learning step
31         self.learn_step_counter = 0
32
33         # initialize zero memory [s, a, r, s_]
34         self.memory = np.zeros((self.memory_size, n_features * 2 + 2))
```

```

37     # consist of [target_net, evaluate_net]
38     self._build_net()
39     t_params = tf.get_collection('target_net_params')
40     e_params = tf.get_collection('eval_net_params')
41     self.replace_target_op = [tf.assign(t, e) for t, e in zip(t_params, e_params)]
42
43     self.sess = tf.Session()
44
45     if output_graph:
46         # $ tensorboard --logdir=logs
47         # tf.train.SummaryWriter soon be deprecated, use following
48         tf.summary.FileWriter("logs/", self.sess.graph)
49
50     self.sess.run(tf.global_variables_initializer())
51     self.cost_his = []
52
53 def _build_net(self):
54     # ----- build evaluate net -----
55     self.s = tf.placeholder(tf.float32, [None, self.n_features], name='s')    # input
56     self.q_target = tf.placeholder(tf.float32, [None, self.n_actions], name='Q_target')    # for calculating loss
57     with tf.variable_scope('eval_net'):
58         # c_names(collections_names) are the collections to store variables
59         c_names, n_l1, w_initializer, b_initializer = \
60             ['eval_net_params', tf.GraphKeys.GLOBAL_VARIABLES], 10, \
61             tf.random_normal_initializer(0., 0.3), tf.constant_initializer(0.1)    # config of layers
62
63     # first layer. collections is used later when assign to target net
64     with tf.variable_scope('l1'):
65         w1 = tf.get_variable('w1', [self.n_features, n_l1], initializer=w_initializer, collections=c_names)
66         b1 = tf.get_variable('b1', [1, n_l1], initializer=b_initializer, collections=c_names)
67         l1 = tf.nn.relu(tf.matmul(self.s, w1) + b1)
68
69     # second layer. collections is used later when assign to target net
70     with tf.variable_scope('l2'):
71         w2 = tf.get_variable('w2', [n_l1, self.n_actions], initializer=w_initializer, collections=c_names)
72         b2 = tf.get_variable('b2', [1, self.n_actions], initializer=b_initializer, collections=c_names)
73         self.q_eval = tf.matmul(l1, w2) + b2
74
75     with tf.variable_scope('loss'):
76         self.loss = tf.reduce_mean(tf.squared_difference(self.q_target, self.q_eval))
77     with tf.variable_scope('train'):
78         self._train_op = tf.train.RMSPropOptimizer(self.lr).minimize(self.loss)
79
80     # ----- build target net -----
81     self.s_ = tf.placeholder(tf.float32, [None, self.n_features], name='s_')    # input
82     with tf.variable_scope('target_net'):
83         # c_names(collections_names) are the collections to store variables
84         c_names = ['target_net_params', tf.GraphKeys.GLOBAL_VARIABLES]
85
86         # first layer. collections is used later when assign to target net
87         with tf.variable_scope('l1'):
88             w1 = tf.get_variable('w1', [self.n_features, n_l1], initializer=w_initializer, collections=c_names)
89             b1 = tf.get_variable('b1', [1, n_l1], initializer=b_initializer, collections=c_names)
90             l1 = tf.nn.relu(tf.matmul(self.s_, w1) + b1)
91
92         # second layer. collections is used later when assign to target net
93         with tf.variable_scope('l2'):
94             w2 = tf.get_variable('w2', [n_l1, self.n_actions], initializer=w_initializer, collections=c_names)
95             b2 = tf.get_variable('b2', [1, self.n_actions], initializer=b_initializer, collections=c_names)
96             self.q_next = tf.matmul(l1, w2) + b2
97
98     def store_transition(self, s, a, r, s_):
99         if not hasattr(self, 'memory_counter'):
100             self.memory_counter = 0
101
102         transition = np.hstack((s, [a, r], s_))
103
104         # replace the old memory with new memory
105         index = self.memory_counter % self.memory_size
106         self.memory[index, :] = transition
107
108         self.memory_counter += 1
109
110     def choose_action(self, observation):
111         # to have batch dimension when feed into tf placeholder
112         observation = observation[np.newaxis, :]
113
114         if np.random.uniform() < self.epsilon:
115             # forward feed the observation and get q value for every actions
116             actions_value = self.sess.run(self.q_eval, feed_dict={self.s: observation})
117             action = np.argmax(actions_value)
118         else:
119             action = np.random.randint(0, self.n_actions)
120
121     return action
122
123     def learn(self):

```

```

123     # check to replace target parameters
124     if self.learn_step_counter % self.replace_target_iter == 0:
125         self.sess.run(self.replace_target_op)
126         print('\'\n\ttarget_params_replaced\'\n')
127
128     # sample batch memory from all memory
129     if self.memory_counter > self.memory_size:
130         sample_index = np.random.choice(self.memory_size, size=self.batch_size)
131     else:
132         sample_index = np.random.choice(self.memory_counter, size=self.batch_size)
133     batch_memory = self.memory[sample_index, :]
134
135     q_next, q_eval = self.sess.run(
136         [self.q_next, self.q_eval],
137         feed_dict={
138             self.s: batch_memory[:, :-self.n_features:], # fixed params
139             self.s: batch_memory[:, :self.n_features], # newest params
140         })
141
142     # change q_target w.r.t q_eval's action
143     q_target = q_eval.copy()
144
145     batch_index = np.arange(self.batch_size, dtype=np.int32)
146     eval_act_index = batch_memory[:, self.n_features].astype(int)
147     reward = batch_memory[:, self.n_features + 1]
148
149     q_target[batch_index, eval_act_index] = reward + self.gamma * np.max(q_next, axis=1)
150
151     # train eval network
152     _, self.cost = self.sess.run([self._train_op, self.loss],
153                                 feed_dict={self.s: batch_memory[:, :self.n_features],
154                                            self.q_target: q_target})
155     self.cost_his.append(self.cost)
156
157     # increasing epsilon
158     self.epsilon = self.epsilon + self.epsilon_increment if self.epsilon < self.epsilon_max else self.epsilon_max
159     self.learn_step_counter += 1

```

学习迭代

```

1 x_threshold = 2.4
2 theta_threshold_radians = 1/15*math.pi
3
4 RL = DeepQNetwork(n_actions=2,
5                 n_features=4,
6                 learning_rate=0.01, e_greedy=0.9,
7                 replace_target_iter=100, memory_size=2000,
8                 e_greedy_increment=0.001,
9
10 total_steps = 0
11
12 for i_episode in range(100):
13     json_req = requests.get(url=url, params={'id': token, 'move': 0}).json()
14     observation = json_req['observation']
15     ep_r = 0
16
17     while True:
18         action = RL.choose_action(np.array(observation))
19
20         json_req = requests.get(url=url, params={'id': token, 'move': action}).json()
21         try: observation_ = json_req['observation']
22         except KeyError:
23             pass
24         print(observation)
25         done = not json_req['status']
26
27         # the smaller theta and closer to center the better
28         x, x_dot, theta, theta_dot = observation_
29         r1 = (x_threshold - abs(x))/x_threshold - 0.8
30         r2 = (theta_threshold_radians - abs(theta))/theta_threshold_radians - 0.5
31         reward = r1 + r2
32
33         RL.store_transition(observation, action, reward, observation_)
34
35         ep_r += reward
36         if total_steps > 1000:
37             RL.learn()
38
39         if done:
40             count = json_req['count']
41             if count == 100:
42                 print(json_req['flag'])
43             else:
44                 print('count:', json_req['count'])
45             break
46
47         observation = observation_
48         total_steps += 1

```

Flag

1 空

- 本文作者: CTFHub
- 本文链接: <https://writeup.ctfhub.com/Challenge/2017/HCTF/Web/wfLYuDnCeIKQNmHXZPqX1.html>
- 版权声明: 本博客所有文章除特别声明外, 均采用 [BY-NC-SA](#) 许可协议。转载请注明出处!

#Challenge #Web #2017 #HCTF

[Deserted place](#)

[Evr_Q](#)