

# E16 Deep Q-Learning (C++/Python)

---

18340215 张天祯

2021 年 1 月 4 日

## Contents

<b>1</b>	<b>Deep Q-Network (DQN)</b>	<b>2</b>
<b>2</b>	<b>Deep Learning Flappy Bird</b>	<b>4</b>
<b>3</b>	<b>Tasks</b>	<b>8</b>
<b>4</b>	<b>Codes and Results</b>	<b>8</b>

# 1 Deep Q-Network (DQN)

We consider tasks in which an agent interacts with an environment  $\mathcal{E}$ , in this case the Atari emulator, in a sequence of actions, observations and rewards. At each time-step the agent selects an action  $a_t$  from the set of legal game actions,  $\mathcal{A} = \{1, \dots, K\}$ . The action is passed to the emulator and modifies its internal state and the game score. In general  $\mathcal{E}$  may be stochastic. The emulator's internal state is not observed by the agent, instead it observes an image  $x_t \in \mathbb{R}^d$  from the emulator, which is a vector of raw pixel values representing the current screen. In addition it receives a reward  $r_t$  representing the change in game score. Note that in general the game score may depend on the whole prior sequence of actions and observations; feedback about an action may only be received after many thousands of time-steps have elapsed.

Since the agent only observes images of the current screen, the task is partially observed and many emulator states are perceptually aliased, i.e. it is impossible to fully understand the current situation from only the current screen  $x_t$ . We therefore consider sequences of actions and observations,  $s_t = x_1, a_1, x_2, \dots, a_{t-1}, x_t$ , and learn game strategies that depend upon these sequences. All sequences in the emulator are assumed to terminate in a finite number of time-steps. This formalism gives rise to a large but finite Markov decision process (MDP) in which each sequence is a distinct state. As a result, we can apply standard reinforcement learning methods for MDPs, simply by using the complete sequence  $s_t$  as the state representation at time  $t$ .

The goal of the agent is to interact with the emulator by selecting actions in a way that maximises future rewards. We make the standard assumption that future rewards are discounted by a factor of  $\gamma$  per time-step, and define the future discounted *return* at time  $t$  as  $R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$ , where  $T$  is the time-step at which the game terminates. We define the optimal action-value function  $Q^*(s, a)$  as the maximum expected return achievable by following any strategy, after seeing some sequence  $s$  and then taking some action  $a$ ,  $Q^*(s, a) = \max_{\pi} \mathbb{E}[R_t | s_t = s, a_t = a, \pi]$ , where  $\pi$  is a policy mapping sequences to actions (or distributions over actions).

The optimal action-value function obeys an important identity known as the *Bellman equation*. This is based on the following intuition: if the optimal value  $Q^*(s', a')$  of the sequence  $s'$  at the next time-step was known for all possible actions  $a'$ , then the optimal strategy is to select the action  $a'$  maximising the expected value of  $r + \gamma Q^*(s', a')$ ,

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}}[r + \gamma \max_{a'} Q^*(s', a') | s, a] \quad (1)$$

The basic idea behind many reinforcement learning algorithms is to estimate the action-value function, by using the Bellman equation as an iterative update,  $Q_{i+1}(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q_i(s', a') | s, a]$ .

Such *value iteration* algorithms converge to the optimal action-value function,  $Q_i \rightarrow Q^*$  as  $i \rightarrow \infty$ . In practice, this basic approach is totally impractical, because the action-value function is estimated separately for each sequence, without any generalisation. Instead, it is common to use a function approximator to estimate the action-value function,  $Q(s, a; \theta) \approx Q^*(s, a)$ . In the reinforcement learning community this is typically a linear function approximator, but sometimes a non-linear function approximator is used instead, such as a neural network. We refer to a neural network function approximator with weights  $\theta$  as a Q-network. A Q-network can be trained by minimising a sequence of loss functions  $L_i(\theta_i)$  that changes at each iteration  $i$ ,

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2], \quad (2)$$

where  $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$  is the target for iteration  $i$  and  $\rho(s, a)$  is a probability distribution over sequences  $s$  and actions  $a$  that we refer to as the *behaviour distribution*. The parameters from the previous iteration  $\theta_{i-1}$  are held fixed when optimising the loss function  $L_i(\theta_i)$ . Note that the targets depend on the network weights; this is in contrast with the targets used for supervised learning, which are fixed before learning begins. Differentiating the loss function with respect to the weights we arrive at the following gradient,

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right]. \quad (3)$$

Rather than computing the full expectations in the above gradient, it is often computationally expedient to optimise the loss function by stochastic gradient descent. If the weights are updated after every time-step, and the expectations are replaced by single samples from the behaviour distribution  $\rho$  and the emulator  $\mathcal{E}$  respectively, then we arrive at the familiar *Q-learning* algorithm.

Note that this algorithm is *model-free*: it solves the reinforcement learning task directly using samples from the emulator  $\mathcal{E}$ , without explicitly constructing an estimate of  $\mathcal{E}$ . It is also *off-policy*: it learns about the greedy strategy  $a = \max_a Q(s, a; \theta)$ , while following a behaviour distribution that ensures adequate exploration of the state space. In practice, the behaviour distribution is often selected by an  $\epsilon$ -greedy strategy that follows the greedy strategy with probability  $1 - \epsilon$  and selects a random action with probability  $\epsilon$ .

---

**Algorithm 1** Deep Q-learning with Experience Replay

---

Initialize replay memory  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

**for** episode = 1,  $M$  **do**

    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$

**for**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$

        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$

        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3

**end for**

**end for**

---

## 2 Deep Learning Flappy Bird

### Overview

This project (<https://github.com/yenchenlin/DeepLearningFlappyBird>) follows the description of the Deep Q Learning algorithm described in *Playing Atari with Deep Reinforcement Learning* and shows that this learning algorithm can be further generalized to the notorious Flappy Bird.

### Installation Dependencies:

- Python 2.7 or 3
- TensorFlow 0.7
- pygame
- OpenCV-Python

### How to Run?

```
git clone https://github.com/yenchenlin1994/DeepLearningFlappyBird.git
```

```
cd DeepLearningFlappyBird
python deep_q_network.py
```

## What is Deep Q-Network?

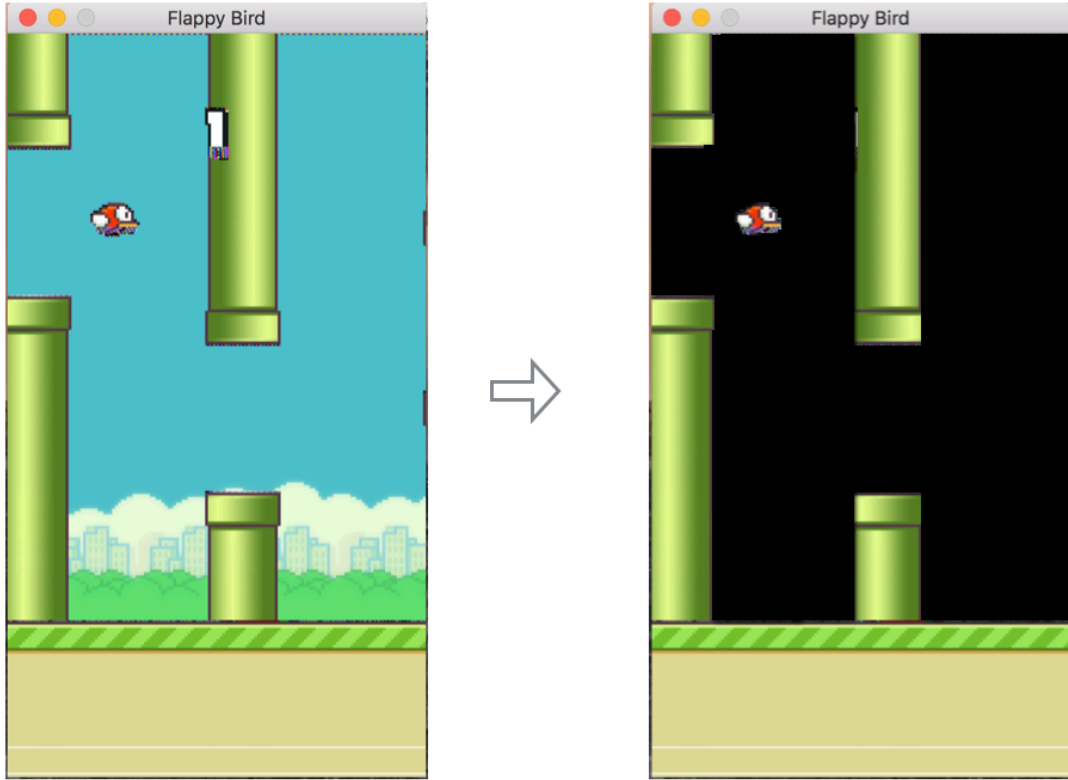
It is a convolutional neural network, trained with a variant of Q-learning, whose input is raw pixels and whose output is a value function estimating future rewards.

For those who are interested in deep reinforcement learning, I highly recommend to read the following post: [Demystifying Deep Reinforcement Learning](#)

## Deep Q-Network Algorithm

The pseudo-code for the Deep Q Learning algorithm can be found below:

```
Initialize replay memory D to size N
Initialize action-value function Q with random weights
for episode = 1, M do
    Initialize state s_1
    for t = 1, T do
        With probability  $\epsilon$  select random action a_t
        otherwise select  $a_t = \arg \max_a Q(s_t, a; \theta_i)$ 
        Execute action a_t in emulator and observe r_t and s_(t+1)
        Store transition (s_t, a_t, r_t, s_(t+1)) in D
        Sample a minibatch of transitions (s_j, a_j, r_j, s_(j+1)) from D
        Set y_j :=
            r_j for terminal s_(j+1)
            r_j +  $\gamma \max_{a'} Q(s_(j+1), a'; \theta_i)$  for non-terminal s_(j+1)
        Perform a gradient step on  $(y_j - Q(s_j, a_j; \theta_i))^2$  with respect to
    end for
end for
```



## Experiments

### Environment

Since deep Q-network is trained on the raw pixel values observed from the game screen at each time step, so removing the background appeared in the original game can make it converge faster.

This process can be visualized as the following figure:

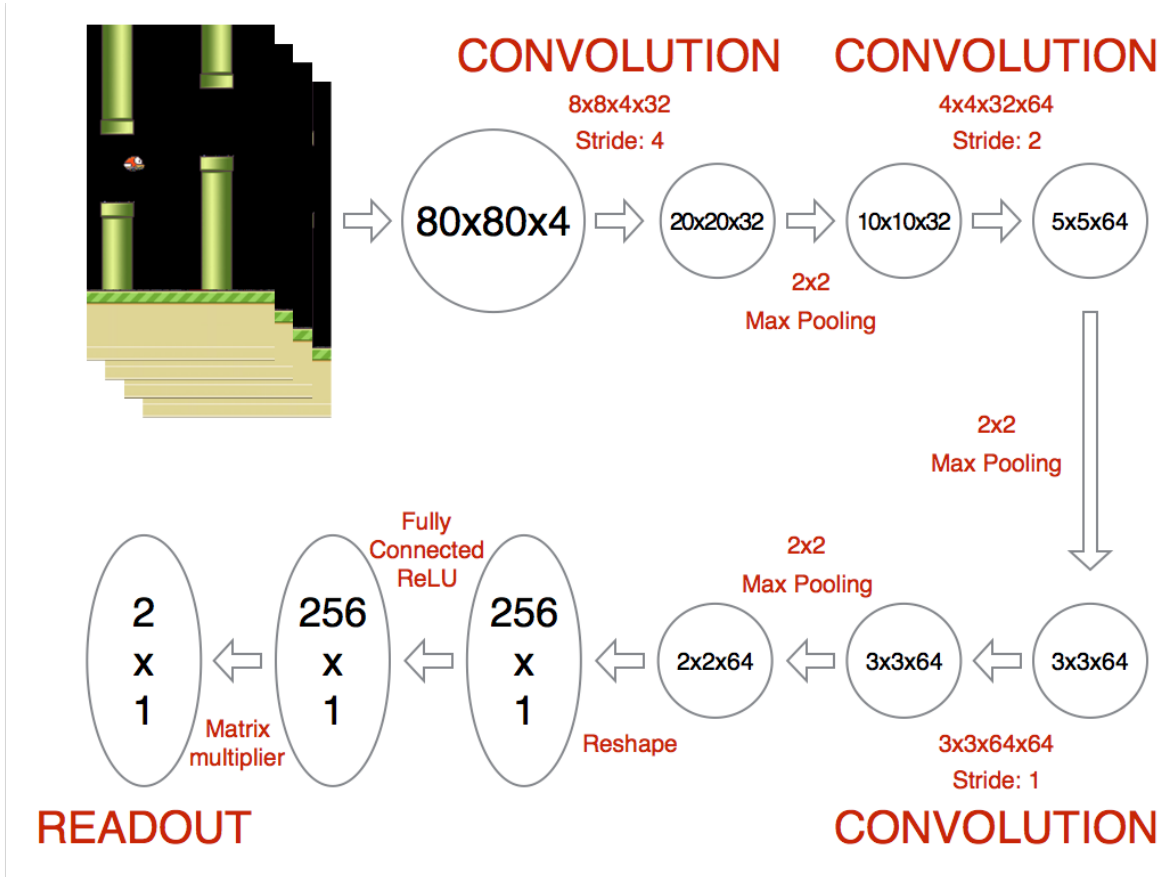
### Network Architecture

I first preprocessed the game screens with following steps:

1. Convert image to grayscale
2. Resize image to 80x80
3. Stack last 4 frames to produce an  $80 \times 80 \times 4$  input array for network

The architecture of the network is shown in the figure below. The first layer convolves the input image with an  $8 \times 8 \times 4 \times 32$  kernel at a stride size of 4. The output is then put through a  $2 \times 2$  max pooling layer. The second layer convolves with a  $4 \times 4 \times 32 \times 64$  kernel at a stride of 2. We then max pool again. The third layer convolves with a  $3 \times 3 \times 64 \times 64$  kernel at a stride of 1. We then max pool one more time. The last hidden layer consists of 256 fully connected ReLU nodes.

The final output layer has the same dimensionality as the number of valid actions which can be



performed in the game, where the 0th index always corresponds to doing nothing. The values at this output layer represent the  $Q$  function given the input state for each valid action. At each time step, the network performs whichever action corresponds to the highest  $Q$  value using a  $\epsilon$  greedy policy.

### Training

At first, I initialize all weight matrices randomly using a normal distribution with a standard deviation of 0.01, then set the replay memory with a max size of 500,00 experiences.

I start training by choosing actions uniformly at random for the first 10,000 time steps, without updating the network weights. This allows the system to populate the replay memory before training begins.

I linearly anneal  $\epsilon$  from 0.1 to 0.0001 over the course of the next 3000,000 frames. The reason why I set it this way is that agent can choose an action every 0.03s (FPS=30) in our game, high  $\epsilon$  will make it **flap** too much and thus keeps itself at the top of the game screen and finally bump the pipe in a clumsy way. This condition will make  $Q$  function converge relatively slow since it only start to look other conditions when  $\epsilon$  is low. However, in other games, initialize  $\epsilon$  to 1 is more reasonable.

During training time, at each time step, the network samples minibatches of size 32 from the

replay memory to train on, and performs a gradient step on the loss function described above using the Adam optimization algorithm with a learning rate of 0.000001. After annealing finishes, the network continues to train indefinitely, with  $\epsilon$  fixed at 0.001.

### 3 Tasks

1. Please implement a DQN to play the Flappy Bird game.
2. You can refer to the codes in <https://github.com/yenchenlin/DeepLearningFlappyBird>
3. Please submit a file named E18\_YourNumber.zip, which should includes the code files and the result pictures, and send it to ai\_2020@foxmail.com

### 4 Codes and Results

```
1  #!/usr/bin/env python
2  from __future__ import print_function
3
4  # import tensorflow as tf
5  import tensorflow.compat.v1 as tf
6  tf.disable_v2_behavior()
7  import cv2
8  import sys
9  sys.path.append("game/")
10 import wrapped_flappy_bird as game
11 import random
12 import numpy as np
13 from collections import deque
14
15 GAME = 'bird' # the name of the game being played for log files
16 ACTIONS = 2 # number of valid actions
17 GAMMA = 0.99 # decay rate of past observations
18 OBSERVE = 100000. # timesteps to observe before training
19 EXPLORE = 2000000. # frames over which to anneal epsilon
20 FINAL_EPSILON = 0.0001 # final value of epsilon
21 INITIAL_EPSILON = 0.0001 # starting value of epsilon
22 REPLAY_MEMORY = 50000 # number of previous transitions to remember
23 BATCH = 32 # size of minibatch
24 FRAME_PER_ACTION = 1
25
26 def weight_variable(shape):
```



```

27     initial = tf.truncated_normal(shape, stddev = 0.01)
28     return tf.Variable(initial)
29
30 def bias_variable(shape):
31     initial = tf.constant(0.01, shape = shape)
32     return tf.Variable(initial)
33
34 def conv2d(x, W, stride):
35     return tf.nn.conv2d(x, W, strides = [1, stride, stride, 1], padding = "SAME")
36
37 def max_pool_2x2(x):
38     return tf.nn.max_pool(x, ksize = [1, 2, 2, 1], strides = [1, 2, 2, 1], padding = "
    SAME")
39
40 def createNetwork():
41     # network weights
42     W_conv1 = weight_variable([8, 8, 4, 32])
43     b_conv1 = bias_variable([32])
44
45     W_conv2 = weight_variable([4, 4, 32, 64])
46     b_conv2 = bias_variable([64])
47
48     W_conv3 = weight_variable([3, 3, 64, 64])
49     b_conv3 = bias_variable([64])
50
51     W_fc1 = weight_variable([1600, 512])
52     b_fc1 = bias_variable([512])
53
54     W_fc2 = weight_variable([512, ACTIONS])
55     b_fc2 = bias_variable([ACTIONS])
56
57     # input layer
58     s = tf.placeholder("float", [None, 80, 80, 4])
59
60     # hidden layers
61     h_conv1 = tf.nn.relu(conv2d(s, W_conv1, 4) + b_conv1)
62     h_pool1 = max_pool_2x2(h_conv1)
63
64     h_conv2 = tf.nn.relu(conv2d(h_pool1, W_conv2, 2) + b_conv2)
65     #h_pool2 = max_pool_2x2(h_conv2)
66

```

```

67     h_conv3 = tf.nn.relu(conv2d(h_conv2, W_conv3, 1) + b_conv3)
68     #h_pool3 = max_pool_2x2(h_conv3)
69
70     #h_pool3_flat = tf.reshape(h_pool3, [-1, 256])
71     h_conv3_flat = tf.reshape(h_conv3, [-1, 1600])
72
73     h_fc1 = tf.nn.relu(tf.matmul(h_conv3_flat, W_fc1) + b_fc1)
74
75     # readout layer
76     readout = tf.matmul(h_fc1, W_fc2) + b_fc2
77
78     return s, readout, h_fc1
79
80 def trainNetwork(s, readout, h_fc1, sess):
81     # define the cost function
82     a = tf.placeholder("float", [None, ACTIONS])
83     y = tf.placeholder("float", [None])
84     readout_action = tf.reduce_sum(tf.multiply(readout, a), reduction_indices=1)
85     cost = tf.reduce_mean(tf.square(y - readout_action))
86     train_step = tf.train.AdamOptimizer(1e-6).minimize(cost)
87
88     # open up a game state to communicate with emulator
89     game_state = game.GameState()
90
91     # store the previous observations in replay memory
92     D = deque()
93
94     # printing
95     a_file = open("logs_" + GAME + "/readout.txt", 'w')
96     h_file = open("logs_" + GAME + "/hidden.txt", 'w')
97
98     # get the first state by doing nothing and preprocess the image to 80x80x4
99     do_nothing = np.zeros(ACTIONS)
100     do_nothing[0] = 1
101     x_t, r_0, terminal = game_state.frame_step(do_nothing)
102     x_t = cv2.cvtColor(cv2.resize(x_t, (80, 80)), cv2.COLOR_BGR2GRAY)
103     ret, x_t = cv2.threshold(x_t, 1, 255, cv2.THRESH_BINARY)
104     s_t = np.stack((x_t, x_t, x_t, x_t), axis=2)
105
106     # saving and loading networks
107     saver = tf.train.Saver()

```

```

108 sess.run(tf.initialize_all_variables())
109 checkpoint = tf.train.get_checkpoint_state("saved_networks")
110 if checkpoint and checkpoint.model_checkpoint_path:
111     saver.restore(sess, checkpoint.model_checkpoint_path)
112     print("Successfully loaded:", checkpoint.model_checkpoint_path)
113 else:
114     print("Could not find old network weights")
115
116 # start training
117 epsilon = INITIAL_EPSILON
118 t = 0
119 while "flappy_bird" != "angry_bird":
120     # choose an action epsilon greedily
121     readout_t = readout.eval(feed_dict={s : [s_t]})[0]
122     a_t = np.zeros([ACTIONS])
123     action_index = 0
124     if t % FRAME_PER_ACTION == 0:
125         if random.random() <= epsilon:
126             print("—————Random Action—————")
127             action_index = random.randrange(ACTIONS)
128             a_t[random.randrange(ACTIONS)] = 1
129         else:
130             action_index = np.argmax(readout_t)
131             a_t[action_index] = 1
132     else:
133         a_t[0] = 1 # do nothing
134
135     # scale down epsilon
136     if epsilon > FINAL_EPSILON and t > OBSERVE:
137         epsilon -= (INITIAL_EPSILON - FINAL_EPSILON) / EXPLORE
138
139     # run the selected action and observe next state and reward
140     x_t1_colored, r_t, terminal = game_state.frame_step(a_t)
141     x_t1 = cv2.cvtColor(cv2.resize(x_t1_colored, (80, 80)), cv2.COLOR_BGR2GRAY)
142     ret, x_t1 = cv2.threshold(x_t1, 1, 255, cv2.THRESH_BINARY)
143     x_t1 = np.reshape(x_t1, (80, 80, 1))
144     #s_t1 = np.append(x_t1, s_t[:, :, 1:], axis = 2)
145     s_t1 = np.append(x_t1, s_t[:, :, :3], axis=2)
146
147     # store the transition in D
148     D.append((s_t, a_t, r_t, s_t1, terminal))

```

```

149     if len(D) > REPLAY_MEMORY:
150         D.popleft()
151
152     # only train if done observing
153     if t > OBSERVE:
154         # sample a minibatch to train on
155         minibatch = random.sample(D, BATCH)
156
157         # get the batch variables
158         s_j_batch = [d[0] for d in minibatch]
159         a_batch = [d[1] for d in minibatch]
160         r_batch = [d[2] for d in minibatch]
161         s_jl_batch = [d[3] for d in minibatch]
162
163         y_batch = []
164         readout_jl_batch = readout.eval(feed_dict = {s : s_jl_batch})
165         for i in range(0, len(minibatch)):
166             terminal = minibatch[i][4]
167             # if terminal, only equals reward
168             if terminal:
169                 y_batch.append(r_batch[i])
170             else:
171                 y_batch.append(r_batch[i] + GAMMA * np.max(readout_jl_batch[i]))
172
173         # perform gradient step
174         train_step.run(feed_dict = {
175             y : y_batch,
176             a : a_batch,
177             s : s_j_batch}
178         )
179
180     # update the old values
181     s_t = s_t1
182     t += 1
183
184     # save progress every 10000 iterations
185     if t % 10000 == 0:
186         saver.save(sess, 'saved_networks/' + GAME + '-dqn', global_step = t)
187
188     # print info
189     state = ""

```

```

190     if t <= OBSERVE:
191         state = "observe"
192     elif t > OBSERVE and t <= OBSERVE + EXPLORE:
193         state = "explore"
194     else:
195         state = "train"
196
197     print("TIMESTEP", t, "/_STATE", state, \
198           "/_EPSILON", epsilon, "/_ACTION", action_index, "/_REWARD", r_t, \
199           "/_Q_MAX_%e" % np.max(readout_t))
200     # write info to files
201     '''
202     if t % 10000 <= 100:
203         a_file.write(",".join([str(x) for x in readout_t]) + '\n')
204         h_file.write(",".join([str(x) for x in h_fc1.eval(feed_dict={s:[s_t]})[0]])
205                        + '\n')
206         cv2.imwrite("logs_tetris/frame" + str(t) + ".png", x_t1)
207     '''
208
209 def playGame():
210     sess = tf.InteractiveSession()
211     s, readout, h_fc1 = createNetwork()
212     trainNetwork(s, readout, h_fc1, sess)
213
214 def main():
215     playGame()
216
217 if __name__ == "__main__":
218     main()

```

本地因为安装了 tensorflow2.\*，而原代码为 tensorflow0.7，所以仅需要在 import 库时做一下额外处理，把库的版本降低即可运行。原代码的结果已经足够好了，没有进行调参。训练次数特别多，结果如下：

```
TIMESTEP 11361 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.275663e+01
TIMESTEP 11362 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.277851e+01
TIMESTEP 11363 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.272854e+01
TIMESTEP 11364 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.275863e+01
TIMESTEP 11365 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.287276e+01
TIMESTEP 11366 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.296755e+01
TIMESTEP 11367 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.284153e+01
TIMESTEP 11368 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.273492e+01
TIMESTEP 11369 / STATE observe / EPSILON 0.0001 / ACTION 1 / REWARD 0.1 / Q_MAX 1.292669e+01
TIMESTEP 11370 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.300509e+01
TIMESTEP 11371 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.303499e+01
TIMESTEP 11372 / STATE observe / EPSILON 0.0001 / ACTION 0 / REWARD 0.1 / Q_MAX 1.297799e+01

Process finished with exit code -1
```

训练时可以看到小鸟飞了近 100 次左右才碰壁，效果拔群。