

Using Tensorflow and Deep Q-Network/Double DQN to Play Breakout

📅 2017-07-09 | 👁

In previous [blog](#), we use the Keras to play the FlappyBird. Similarity, we will use another deep learning toolkit Tensorflow to develop the DQN and Double DQN and to play the another game Breakout (Atari 3600).

Here, we will use the OpenAI gym toolkit to construct out environment. Detail implementation is as follows:

```
1 env = gym.envs.make("Breakout-v0")
```

And then we look some demos:

```
1 print("Action space size: {}".format(env.action_space.n))
2 # print(env.get_action_meanings())
3
4 observation = env.reset()
5 print("Observation space shape: {}".format(observation.shape))
6
7 plt.figure()
8 plt.imshow(env.render(mode='rgb_array'))
9
10 [env.step(2) for x in range(1)]
11 plt.figure()
12 plt.imshow(env.render(mode='rgb_array'))
13
14 env.render(close=True)
```

 breakout-env

For deep learning purpose, we need to crop the image to a square image:

 # Check out what a cropped image looks like

```
2 plt.imshow(observation[34:-16,:,:])
```

cropped-breakout-image

Not bad !

Ok, now let us to use the Tensorflow to develop the DQN algorithm first.

First of all, we need to reference some packages and initialize the environment.

```
1  %matplotlib inline
2
3  import gym
4  from gym.wrappers import Monitor
5  import itertools
6  import numpy as np
7  import os
8  import random
9  import sys
10 import tensorflow as tf
11
12 if "../" not in sys.path:
13     sys.path.append("../")
14
15 from lib import plotting
16 from collections import deque, namedtuple
17
18 env = gym.envs.make("Breakout-v0")
19 # Atari Actions: 0 (noop), 1 (fire), 2 (left) and 3 (right) are valid actions
20 VALID_ACTIONS = [0, 1, 2, 3]
```

As mentioned above, we need to crop the image and preprocess the input before feed the raw image into the algorithm. So we define a **StateProcessor** class to do this.

```
1  class StateProcessor():
2      """
3      Processes a raw Atari images. Resizes it and converts it to grayscale.
4      """
5      def __init__(self):
6          # Build the Tensorflow graph
7          with tf.variable_scope("state_processor"):
8              self.input_state = tf.placeholder(shape=[210, 160, 3], dtype=tf.uint8)
9              self.output = tf.image.rgb_to_grayscale(self.input_state)
10             self.output = tf.image.crop_to_bounding_box(self.output, 34, 0, 160,
```

```

11         self.output = tf.image.resize_images(
12             self.output, [84, 84], method=tf.image.ResizeMethod.NEAREST_NEI(
13             self.output = tf.squeeze(self.output)
14
15     def process(self, sess, state):
16         """
17         Args:
18             sess: A Tensorflow session object
19             state: A [210, 160, 3] Atari RGB State
20
21         Returns:
22             A processed [84, 84, 1] state representing grayscale values.
23         """
24         return sess.run(self.output, { self.input_state: state })

```

We first convert the image to gray image and then resize it to 84 by 84 (the size DQN paper used). Then, we construct the neural network to estimate the value function. The structure of the network as the same as the DQN paper.

```

1  class Estimator():
2      """Q-Value Estimator neural network.
3
4      This network is used for both the Q-Network and the Target Network.
5      """
6
7      def __init__(self, scope="estimator", summaries_dir=None):
8          self.scope = scope
9          # Writes Tensorboard summaries to disk
10         self.summary_writer = None
11         with tf.variable_scope(scope):
12             # Build the graph
13             self._build_model()
14             if summaries_dir:
15                 summary_dir = os.path.join(summaries_dir, "summaries_{}".format(
16                     if not os.path.exists(summary_dir):
17                         os.makedirs(summary_dir)
18                 self.summary_writer = tf.summary.FileWriter(summary_dir)
19
20     def _build_model(self):
21         """
22         Builds the Tensorflow graph.
23         """
24
25         # Placeholders for our input
26         # Our input are 4 RGB frames of shape 160, 160 each
27         self.X_pl = tf.placeholder(shape=[None, 84, 84, 4], dtype=tf.uint8, name=
28         # The TD target value
29         self.y_pl = tf.placeholder(shape=[None], dtype=tf.float32, name="y")

```

```

30     # Integer id of which action was selected
31     self.actions_pl = tf.placeholder(shape=[None], dtype=tf.int32, name="actions")
32
33     X = tf.to_float(self.X_pl) / 255.0
34     batch_size = tf.shape(self.X_pl)[0]
35
36     # Three convolutional layers
37     conv1 = tf.contrib.layers.conv2d(
38         X, 32, 8, 4, activation_fn=tf.nn.relu)
39     conv2 = tf.contrib.layers.conv2d(
40         conv1, 64, 4, 2, activation_fn=tf.nn.relu)
41     conv3 = tf.contrib.layers.conv2d(
42         conv2, 64, 3, 1, activation_fn=tf.nn.relu)
43
44     # Fully connected layers
45     flattened = tf.contrib.layers.flatten(conv3)
46     fc1 = tf.contrib.layers.fully_connected(flattened, 512)
47     self.predictions = tf.contrib.layers.fully_connected(fc1, len(VALID_ACTIONS))
48
49     # Get the predictions for the chosen actions only
50     gather_indices = tf.range(batch_size) * tf.shape(self.predictions)[1] + self.actions_pl
51     self.action_predictions = tf.gather(tf.reshape(self.predictions, [-1]), gather_indices)
52
53     # Calculate the loss
54     self.losses = tf.squared_difference(self.y_pl, self.action_predictions)
55     self.loss = tf.reduce_mean(self.losses)
56
57     # Optimizer Parameters from original paper
58     self.optimizer = tf.train.RMSPropOptimizer(0.00025, 0.99, 0.0, 1e-6)
59     self.train_op = self.optimizer.minimize(self.loss, global_step=tf.contrib.framework.get_global_step())
60
61     # Summaries for Tensorboard
62     self.summaries = tf.summary.merge([
63         tf.summary.scalar("loss", self.loss),
64         tf.summary.histogram("loss_hist", self.losses),
65         tf.summary.histogram("q_values_hist", self.predictions),
66         tf.summary.scalar("max_q_value", tf.reduce_max(self.predictions))
67     ])
68
69     def predict(self, sess, s):
70         """
71         Predicts action values.
72
73         Args:
74             sess: Tensorflow session
75             s: State input of shape [batch_size, 4, 160, 160, 3]
76
77         Returns:
78             Tensor of shape [batch_size, NUM_VALID_ACTIONS] containing the estimated
79             action values.
80         """
81         return sess.run(self.predictions, { self.X_pl: s })

```

```

82
83     def update(self, sess, s, a, y):
84         """
85         Updates the estimator towards the given targets.
86
87         Args:
88             sess: Tensorflow session object
89             s: State input of shape [batch_size, 4, 160, 160, 3]
90             a: Chosen actions of shape [batch_size]
91             y: Targets of shape [batch_size]
92
93         Returns:
94             The calculated loss on the batch.
95         """
96         feed_dict = { self.X_pl: s, self.y_pl: y, self.actions_pl: a }
97         summaries, global_step, _, loss = sess.run(
98             [self.summaries, tf.contrib.framework.get_global_step(), self.train
99             feed_dict]
100         if self.summary_writer:
101             self.summary_writer.add_summary(summaries, global_step)
102         return loss

```

As mentioned in DQN paper, there are two network that share the same parameters in DQN algorithm. We need to copy the parameters to the target network on each t steps.

```

1  def copy_model_parameters(sess, estimator1, estimator2):
2      """
3      Copies the model parameters of one estimator to another.
4
5      Args:
6          sess: Tensorflow session instance
7          estimator1: Estimator to copy the paramters from
8          estimator2: Estimator to copy the parameters to
9      """
10     e1_params = [t for t in tf.trainable_variables() if t.name.startswith(estimator1)]
11     e1_params = sorted(e1_params, key=lambda v: v.name)
12     e2_params = [t for t in tf.trainable_variables() if t.name.startswith(estimator2)]
13     e2_params = sorted(e2_params, key=lambda v: v.name)
14
15     update_ops = []
16     for e1_v, e2_v in zip(e1_params, e2_params):
17         op = e2_v.assign(e1_v)
18         update_ops.append(op)
19
20     sess.run(update_ops)

```

We also need a policy to take an action.

```
1 def make_epsilon_greedy_policy(estimator, nA):
2     """
3     Creates an epsilon-greedy policy based on a given Q-function approximator and
4     an action selector.
5     Args:
6         estimator: An estimator that returns q values for a given state
7         nA: Number of actions in the environment.
8
9     Returns:
10        A function that takes the (sess, observation, epsilon) as an argument and
11        returns the probabilities for each action in the form of a numpy array of length
12        nA.
13    """
14    def policy_fn(sess, observation, epsilon):
15        A = np.ones(nA, dtype=float) * epsilon / nA
16        q_values = estimator.predict(sess, np.expand_dims(observation, 0))[0]
17        best_action = np.argmax(q_values)
18        A[best_action] += (1.0 - epsilon)
19        return A
20    return policy_fn
```

Now let us to develop the DQN algorithm (we skip the details here because we explained it earlier).

```
1 def deep_q_learning(sess,
2                     env,
3                     q_estimator,
4                     target_estimator,
5                     state_processor,
6                     num_episodes,
7                     experiment_dir,
8                     replay_memory_size=500000,
9                     replay_memory_init_size=50000,
10                    update_target_estimator_every=10000,
11                    discount_factor=0.99,
12                    epsilon_start=1.0,
13                    epsilon_end=0.1,
14                    epsilon_decay_steps=500000,
15                    batch_size=32,
16                    record_video_every=50):
17     """
18     Q-Learning algorithm for fff-policy TD control using Function Approximation.
19     Finds the optimal greedy policy while following an epsilon-greedy policy.
20
21     Args:
22         sess: Tensorflow Session object
```

```

23     env: OpenAI environment
24     q_estimator: Estimator object used for the q values
25     target_estimator: Estimator object used for the targets
26     state_processor: A StateProcessor object
27     num_episodes: Number of episodes to run for
28     experiment_dir: Directory to save Tensorflow summaries in
29     replay_memory_size: Size of the replay memory
30     replay_memory_init_size: Number of random experiences to sampel when in
31         the replay memory.
32     update_target_estimator_every: Copy parameters from the Q estimator to
33         target estimator every N steps
34     discount_factor: Lambda time discount factor
35     epsilon_start: Chance to sample a random action when taking an action.
36         Epsilon is decayed over time and this is the start value
37     epsilon_end: The final minimum value of epsilon after decaying is done
38     epsilon_decay_steps: Number of steps to decay epsilon over
39     batch_size: Size of batches to sample from the replay memory
40     record_video_every: Record a video every N episodes
41
42 Returns:
43     An EpisodeStats object with two numpy arrays for episode_lengths and episode_rewards
44     """
45
46     Transition = namedtuple("Transition", ["state", "action", "reward", "next_state"])
47
48     # The replay memory
49     replay_memory = []
50
51     # Keeps track of useful statistics
52     stats = plotting.EpisodeStats(
53         episode_lengths=np.zeros(num_episodes),
54         episode_rewards=np.zeros(num_episodes))
55
56     # Create directories for checkpoints and summaries
57     checkpoint_dir = os.path.join(experiment_dir, "checkpoints")
58     checkpoint_path = os.path.join(checkpoint_dir, "model")
59     monitor_path = os.path.join(experiment_dir, "monitor")
60
61     if not os.path.exists(checkpoint_dir):
62         os.makedirs(checkpoint_dir)
63     if not os.path.exists(monitor_path):
64         os.makedirs(monitor_path)
65
66     saver = tf.train.Saver()
67     # Load a previous checkpoint if we find one
68     latest_checkpoint = tf.train.latest_checkpoint(checkpoint_dir)
69     if latest_checkpoint:
70         print("Loading model checkpoint {}...\n".format(latest_checkpoint))
71         saver.restore(sess, latest_checkpoint)
72
73     # Get the current time step
74     total_t = sess.run(tf.contrib.framework.get_global_step())

```

```

75
76     # The epsilon decay schedule
77     epsilons = np.linspace(epsilon_start, epsilon_end, epsilon_decay_steps)
78
79     # The policy we're following
80     policy = make_epsilon_greedy_policy(
81         q_estimator,
82         len(VALID_ACTIONS))
83
84     # Populate the replay memory with initial experience
85     print("Populating replay memory...")
86     state = env.reset()
87     state = state_processor.process(sess, state)
88     state = np.stack([state] * 4, axis=2)
89     for i in range(replay_memory_init_size):
90         action_probs = policy(sess, state, epsilons[min(total_t, epsilon_decay_
91             action = np.random.choice(np.arange(len(action_probs)), p=action_probs)
92         next_state, reward, done, _ = env.step(VALID_ACTIONS[action])
93         next_state = state_processor.process(sess, next_state)
94         next_state = np.append(state[:, :, 1:], np.expand_dims(next_state, 2), a
95         replay_memory.append(Transition(state, action, reward, next_state, done
96         if done:
97             state = env.reset()
98             state = state_processor.process(sess, state)
99             state = np.stack([state] * 4, axis=2)
100         else:
101             state = next_state
102
103     # Record videos
104     # Add env Monitor wrapper
105     env = Monitor(env, directory=monitor_path, video_callable=lambda count: co
106
107     for i_episode in range(num_episodes):
108
109         # Save the current checkpoint
110         saver.save(tf.get_default_session(), checkpoint_path)
111
112         # Reset the environment
113         state = env.reset()
114         state = state_processor.process(sess, state)
115         state = np.stack([state] * 4, axis=2)
116         loss = None
117
118         # One step in the environment
119         for t in itertools.count():
120
121             # Epsilon for this time step
122             epsilon = epsilons[min(total_t, epsilon_decay_steps-1)]
123
124             # Add epsilon to Tensorboard
125             episode_summary = tf.Summary()
126

```



```

127 episode_summary.value.add(simple_value=epsilon, tag="epsilon")
128 q_estimator.summary_writer.add_summary(episode_summary, total_t)
129
130 # Maybe update the target estimator
131 if total_t % update_target_estimator_every == 0:
132     copy_model_parameters(sess, q_estimator, target_estimator)
133     print("\nCopied model parameters to target network.")
134
135 # Print out which step we're on, useful for debugging.
136 print("\rStep {} ({} ) @ Episode {}/{}", loss: {}".format(
137     t, total_t, i_episode + 1, num_episodes, loss), end="")
138 sys.stdout.flush()
139
140 # Take a step
141 action_probs = policy(sess, state, epsilon)
142 action = np.random.choice(np.arange(len(action_probs)), p=action_probs)
143 next_state, reward, done, _ = env.step(VALID_ACTIONS[action])
144 next_state = state_processor.process(sess, next_state)
145 next_state = np.append(state[:, :, 1:], np.expand_dims(next_state, 2), axis=-1)
146
147 # If our replay memory is full, pop the first element
148 if len(replay_memory) == replay_memory_size:
149     replay_memory.pop(0)
150
151 # Save transition to replay memory
152 replay_memory.append(Transition(state, action, reward, next_state, done))
153
154 # Update statistics
155 stats.episode_rewards[i_episode] += reward
156 stats.episode_lengths[i_episode] = t
157
158 # Sample a minibatch from the replay memory
159 samples = random.sample(replay_memory, batch_size)
160 states_batch, action_batch, reward_batch, next_states_batch, done_batch = zip(*samples)
161
162 # Calculate q values and targets
163 q_values_next = target_estimator.predict(sess, next_states_batch)
164 targets_batch = reward_batch + np.invert(done_batch).astype(np.float32) * q_values_next
165
166 # Perform gradient descent update
167 states_batch = np.array(states_batch)
168 loss = q_estimator.update(sess, states_batch, action_batch, targets_batch)
169
170 if done:
171     break
172
173 state = next_state
174 total_t += 1
175
176 # Add summaries to tensorboard
177 episode_summary = tf.Summary()
178 episode_summary.value.add(simple_value=stats.episode_rewards[i_episode], tag="episode_rewards")

```

```

179     episode_summary.value.add(simple_value=stats.episode_lengths[i_episode])
180     q_estimator.summary_writer.add_summary(episode_summary, total_t)
181     q_estimator.summary_writer.flush()
182
183     yield total_t, plotting.EpisodeStats(
184         episode_lengths=stats.episode_lengths[:i_episode+1],
185         episode_rewards=stats.episode_rewards[:i_episode+1])
186
187     return stats

```

Finally, run it.

```

1  tf.reset_default_graph()
2
3  # Where we save our checkpoints and graphs
4  experiment_dir = os.path.abspath("./experiments/{}".format(env.spec.id))
5
6  # Create a global step variable
7  global_step = tf.Variable(0, name='global_step', trainable=False)
8
9  # Create estimators
10 q_estimator = Estimator(scope="q", summaries_dir=experiment_dir)
11 target_estimator = Estimator(scope="target_q")
12
13 # State processor
14 state_processor = StateProcessor()
15
16 # Run it!
17 with tf.Session() as sess:
18     sess.run(tf.global_variables_initializer())
19     for t, stats in deep_q_learning(sess,
20                                     env,
21                                     q_estimator=q_estimator,
22                                     target_estimator=target_estimator,
23                                     state_processor=state_processor,
24                                     experiment_dir=experiment_dir,
25                                     num_episodes=10000,
26                                     replay_memory_size=500000,
27                                     replay_memory_init_size=50000,
28                                     update_target_estimator_every=10000,
29                                     epsilon_start=1.0,
30                                     epsilon_end=0.1,
31                                     epsilon_decay_steps=500000,
32                                     discount_factor=0.99,
33                                     batch_size=32):
34
35         print("\nEpisode Reward: {}".format(stats.episode_rewards[-1]))

```

Next, we will develop the Double-DQN algorithm. This algorithm only has a little changes.

In DQN `q_learning` method,

```
1  # Calculate q values and targets
2  q_values_next = target_estimator.predict(sess, next_states_batch)
3  targets_batch = reward_batch + np.invert(done_batch).astype(np.float32) * discou
```

we just change these codes to,

```
1  # Calculate q values and targets
2  # This is where Double Q-Learning comes in!
3  q_values_next = q_estimator.predict(sess, next_states_batch)
4  best_actions = np.argmax(q_values_next, axis=1)
5  q_values_next_target = target_estimator.predict(sess, next_states_batch)
6  targets_batch = reward_batch + np.invert(done_batch).astype(np.float32) * discou
```

deep learning # machine learning # reinforcement learning # deep reinforcement learning

◀ Summary of Papers

Policy Gradient Methods ▶

© 2018 ♥ Ewan Li

Powered by [Hexo](#) | Theme - [NexT.Mist](#)

👤 本站访客数 人次 | 👁 本站总访问量 次

