#ELEN 6885 Reinforcement Learning Coding Assignment (Part 4)# There are a lot of official and unofficial tutorials about Tensorflow, and there are also many open-source projects written in Tensorflow. You can refer to those resources according to your interest. In this part of homework 4, only knowledge of Deep Reinforcement Learning and basic programming skills will be needed.

Please put your code into the block marked by\ ############################\ # YOUR CODE STARTS HERE\ # YOUR CODE ENDS HERE\ ###########################\ Normally you don't need to edit anything outside of the block. If you do want to edit something, please use a similar manner to mark you edits.

In [0]:

```python
import numpy as np
import tensorflow as tf

# DQN
class DQN:
    def __init__(
        self,
        actions_num,
        state_size,
        learning_rate = 0.001,
        gamma = 0.99,
        epsilon_min = 0.05,
        epsilon_start = 0.9,
        replace_target_iter = 300,
        memory_size = 500,
        batch_size = 2,
        epsilon_increment = None,
    ):
        self.actions_num = actions_num
        self.state_size = state_size
        self.lr = learning_rate
        self.gamma = gamma
        self.epsilon_min = epsilon_min
        self.replace_target_iter = replace_target_iter
        self.memory_size = memory_size
        self.batch_size = batch_size
        self.epsilon_increment = epsilon_increment
        self.epsilon = epsilon_start if epsilon_increment is not None e
        self.save_model_path = './weights/DQN_model.ckpt'
        self.memory_counter = 0

        # learned steps counter
        self.steps_counter = 0

        # initialize memory [s, a, r, s_, done]
        self.memory = np.zeros((self.memory_size, state_size * 2 + 3))

        # build target_net and q_net
        self.build_net()
        t_params = tf.get_collection('target_net_params')
        q_params = tf.get_collection('q_net_params')
        self.replace_target = [tf.assign(t, q) for t, q in zip(t_params

        # gpu setting
        config = tf.ConfigProto(log_device_placement=False, allow_soft_
        config.gpu_options.per_process_gpu_memory_fraction = 0.6
        self.sess = tf.Session(config=config)

        self.sess.run(tf.global_variables_initializer())

    def build_net(self):
        # build q_net
        self.state = tf.placeholder(tf.float32, [None, self.state_size],
        self.q_target = tf.placeholder(tf.float32, [None, self.actions_nu
        with tf.variable_scope('q_net'):
            # c_names(collections_names) are the collections to store varia
```

```
57        c_names, neurons_layer_1, w_initializer, b_initializer = \
58          ['q_net_params', tf.GraphKeys.GLOBAL_VARIABLES], 100, \
59          tf.random_normal_initializer(0., 0.3), tf.constant_initialize
60
61        # layer 1
62        with tf.variable_scope('layer_1'):
63          w_layer_1 = tf.get_variable('w_layer_1', [self.state_size, ne
64          b_layer_1 = tf.get_variable('b_layer_1', [1, neurons_layer_1]
65          layer_1 = tf.nn.relu(tf.matmul(self.state, w_layer_1) + b_lay
66
67        # layer 2
68        with tf.variable_scope('layer_2'):
69          w_layer_2 = tf.get_variable('w_layer_2', [neurons_layer_1, se
70          b_layer_2 = tf.get_variable('b_layer_2', [1, self.actions_num
71          self.q_value = tf.matmul(layer_1, w_layer_2) + b_layer_2
72
73      with tf.variable_scope('loss'):
74        self.loss = tf.reduce_mean(tf.squared_difference(self.q_target,
75      with tf.variable_scope('train'):
76        self._train_op = tf.train.AdamOptimizer(self.lr).minimize(self.
77
78      # build target_net
79      self.state_t = tf.placeholder(tf.float32, [None, self.state_size]
80      with tf.variable_scope('target_net'):
81        # c_names(collections_names) are the collections to store varia
82        c_names = ['target_net_params', tf.GraphKeys.GLOBAL_VARIABLES]
83
84        # layer 1
85        with tf.variable_scope('layer_1'):
86          w_layer_1 = tf.get_variable('w_layer_1', [self.state_size, ne
87          b_layer_1 = tf.get_variable('b_layer_1', [1, neurons_layer_1]
88          layer_1 = tf.nn.relu(tf.matmul(self.state_t, w_layer_1) + b_l
89
90        # layer 2
91
92        ###########################
93        # YOUR CODE STARTS HERE
94        with tf.variable_scope('layer_2'):
95          w_layer_2 = tf.get_variable('w_layer_2', [neurons_layer_1, se
96          b_layer_2 = tf.get_variable('b_layer_2', [1, self.actions_num
97          self.q_next = tf.matmul(layer_1, w_layer_2) + b_layer_2
98
99        # YOUR CODE ENDS HERE
100       ###########################
101
102   def store_transition(self, s, a, r, s_, done):
103     s=s.reshape(-1)
104     s_=s_.reshape(-1)
105     transition = np.hstack((s, [a, r], s_, done))
106     # replace the old memory with new observations
107     index = self.memory_counter % self.memory_size
108     self.memory[index, :] = transition
109
110     self.memory_counter += 1
111
112   def choose_action(self, observation):
113     # to have batch dimension when fed into tf placeholder
```

```python
114         observation = observation[np.newaxis, :]
115       # epsilon-greedy
116       if np.random.uniform() > self.epsilon:
117         action_values = self.sess.run(self.q_value, feed_dict={self.sta
118         action = np.argmax(action_values)
119       else:
120         action = np.random.randint(0, self.actions_num)
121       return action
122
123   def learn(self):
124     # replace target parameters every once a while
125     if self.steps_counter % self.replace_target_iter == 0:
126       self.sess.run(self.replace_target)
127
128     # sample a batch from the memory
129     if self.memory_counter > self.memory_size:
130       sample_index = np.random.choice(self.memory_size, size=self.bat
131     else:
132       sample_index = np.random.choice(self.memory_counter, size=self.
133     batch_memory = self.memory[sample_index, :]
134
135     q_next, q_value = self.sess.run(
136       [self.q_next, self.q_value],
137       feed_dict={
138         self.state_t: batch_memory[:, -self.state_size-1:-1],  # fixe
139         self.state: batch_memory[:, :self.state_size],  # newest para
140       })
141
142     # calculate q_target
143     q_target = q_value.copy()
144
145
146     # only change the action-values of this batch, because we only ca
147     batch_index = np.arange(self.batch_size, dtype=np.int32)
148     act_index = batch_memory[:, self.state_size].astype(int)
149     reward = batch_memory[:, self.state_size + 1]
150     done = batch_memory[:, -1]
151     ###########################
152     # YOUR CODE STARTS HERE
153     q_target[batch_index, act_index] = reward + self.gamma * np.max(q
154     for i in range(done.shape[0]):
155       if done[i] == 1.0:
156         q_target[i,:] = reward[i]
157
158     # YOUR CODE ENDS HERE
159     ###########################
160
161     # train q_net
162     _, self.cost = self.sess.run([self._train_op, self.loss],
163                                   feed_dict={self.state: batch_memory
164                                              self.q_target: q_target})
165     # change epsilon
166     self.epsilon = self.epsilon - self.epsilon_increment if self.epsi
167     self.steps_counter += 1
168
169   def store(self):
170     saver = tf.train.Saver()
```

```
171        saver.save(self.sess, self.save_model_path)
172
173  ▾    def restore(self):
174          saver = tf.train.Saver()
175          saver.restore(self.sess, self.save_model_path)
176
177
178
179
180
```

Type *Markdown* and LaTeX: $\alpha^2$

In [0]:
```
 1  import gym
 2  # cart pole gym environment
 3  env = gym.make("CartPole-v0")
 4  env._max_episode_steps = 500
 5  # state and action space
 6  print(env.action_space)
 7  print(env.observation_space)
 8  # observation
 9  env.reset()
10  # state, reward, done, info
11  print(env.step(1))
```

```
Discrete(2)
Box(4,)
(array([ 0.03984107,  0.17240988, -0.04123799, -0.2650715 ]), 1.0, False,
{})
```

In [0]:
```python
# play the game and train the network
np.set_printoptions(threshold=np.inf)
episode_length_set = []
tf.reset_default_graph()
total_time_steps = 100000

RL = DQN(actions_num = 2, gamma = 0.99,
         state_size = 4, epsilon_start = 1,
         learning_rate = 1e-3, epsilon_min = 0.01,
         replace_target_iter = 100, memory_size = 5000,
         epsilon_increment = 0.00001,)

new_state = env.reset()
done = False
episode_length_counter = 0
for step in range(total_time_steps):
    ###########################
    # YOUR CODE STARTS HERE
    action = RL.choose_action(new_state)
    next_state, reward, done,_ = env.step(action)
    RL.store_transition(new_state, action, reward, next_state, done)
    if done:
        new_state = env.reset()
        episode_length_set.append(episode_length_counter)
        episode_length_counter =0
    else:
        new_state = next_state

    # YOUR CODE ENDS HERE
    ###########################
    #print(step)
    if step > 200:
      RL.learn()
    episode_length_counter += 1
    if episode_length_counter == 500:
      print('we hit 500')
      RL.store()
```
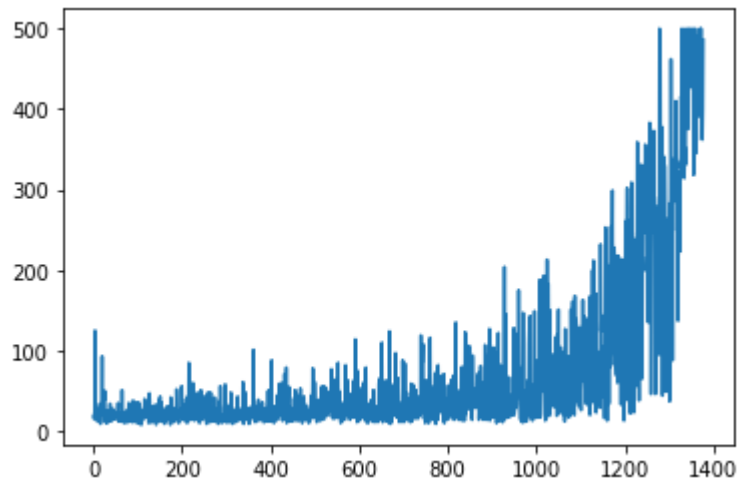
```
we hit 500
we hit 500
we hit 500
we hit 500
we hit 500
we hit 500
we hit 500
we hit 500
we hit 500
```

In [0]:
```python
from matplotlib import pyplot as plt
plt.plot(episode_length_set)
```

Out[146]: [<matplotlib.lines.Line2D at 0x7f5fc28c7978>]

```python
In [0]:   1 ▾  # test our network
          2    tf.reset_default_graph()
          3 ▾  RL = DQN(actions_num = 2, gamma = 1,
          4             state_size = 4, epsilon_start = 1,
          5             learning_rate = 1e-3, epsilon_min = 0,
          6             replace_target_iter = 100, memory_size = 5000,
          7             epsilon_increment = None,)
          8    # load saved parameters
          9    RL.restore()
         10    # run 100 trails and print how long can the agent hold the cart pole f
         11    length_counter = []
         12 ▾  for i in range(100):
         13       ###########################
         14       # YOUR CODE STARTS HERE
         15       new_state = env.reset()
         16       episode_length = 0
         17       done = False
         18 ▾     while not done:
         19          action = RL.choose_action(new_state)
         20          next_state, reward, done,_ = env.step(action)
         21          RL.store_transition(new_state, action, reward, next_state, done)
         22          episode_length += 1
         23 ▾        if done:
         24             new_state = env.reset()
         25             length_counter.append(episode_length)
         26 ▾        else:
         27             new_state = next_state
         28
         29       # YOUR CODE ENDS HERE
         30       ###########################
         31
```
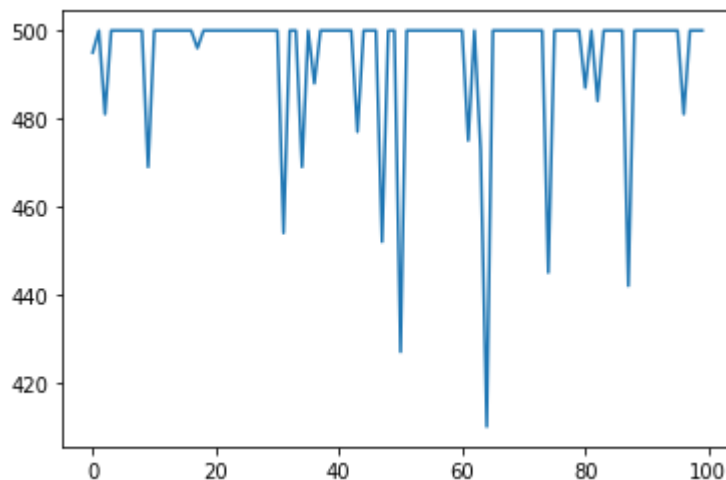
INFO:tensorflow:Restoring parameters from ./weights/DQN_model.ckpt

```python
In [0]:   1    plt.plot(length_counter)
```

Out[148]:  [<matplotlib.lines.Line2D at 0x7f5fc1533a20>]

In [0]:    1    `print(length_counter)`

```
[495, 500, 481, 500, 500, 500, 500, 500, 500, 469, 500, 500, 500, 500, 50
0, 500, 500, 496, 500, 500, 500, 500, 500, 500, 500, 500, 500, 500, 500,
500, 500, 454, 500, 500, 469, 500, 488, 500, 500, 500, 500, 500, 500, 47
7, 500, 500, 500, 452, 500, 500, 427, 500, 500, 500, 500, 500, 500, 500,
500, 500, 500, 475, 500, 473, 410, 500, 500, 500, 500, 500, 500, 500, 50
0, 500, 445, 500, 500, 500, 500, 500, 487, 500, 484, 500, 500, 500, 500,
442, 500, 500, 500, 500, 500, 500, 500, 500, 481, 500, 500, 500]
```

You may find that the episode length doesn't stably improve as more training time is given. You can read chapter 3.2 of this paper https://arxiv.org/pdf/1711.07478.pdf (https://arxiv.org/pdf/1711.07478.pdf) if you are interested.

In [0]:    1    `print(length_counter)`