

CSC411 - Project #4

Yui Chit (Michael) Wong - 999806232
Yijin (Catherine) Wang - 998350476

March 27, 2017

Part 1

Question Explain precisely why the code corresponds to the pseudocode below. Specifically, in your report, explain how all the terms (G_t , π , and the update to θ) are computed, quoting the relevant lines of Python.

REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

```
Input: a differentiable policy parameterization  $\pi(a|s, \theta), \forall a \in \mathcal{A}, s \in \mathcal{S}, \theta \in \mathbb{R}^n$ 
Initialize policy weights  $\theta$ 
Repeat forever:
    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$ 
    For each step of the episode  $t = 0, \dots, T - 1$ :
         $G_t \leftarrow$  return from step  $t$ 
         $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t|S_t, \theta)$ 
```

Pseudocode

Answer As mentioned on the assignment page, the policy function π_{θ} is implemented with a single-hidden-layer of neural network. Since the actions for the bipedal walker is continuous, we have to use a Gaussian distribution on π . Thus we pass the hidden layer into two separately fully connected output, which represents the μ and σ to the normal distribution. The activation function are *tanh* and *softplus* (variation on *ReLU*) respectively. The *sigma* value are also clipped if it is too small or big.

As for the weight initialization, if there is no weight saved from the previous run, we will initialize the weight θ . There are one w and b for each of the layers (hidden, μ , σ). When initializing the weight to each layer, the program uses *xavierinitialization*, another variation of random weight initialization that keep the scale of the gradients in roughly the same scale.

Once everything is initialized, we will start training. For each iteration, we will reset the environment (line 118), then generate the states, actions, and rewards from time 0 to time T . When generating the actions, we will randomly sample a π_{sample} from the π normal distribution (line 133-135). Then based on the π_{sample} , we will generate the corresponding action, and using the action, the new state, reward would be generated. We will keep track of all the states, actions, and rewards in 3 lists (*ep_states*, *ep_actions*, and *ep_rewards*). We will also keep track of the total discounted rewards using the variable G (line 137). Then to obtain G_t , the discounted reward starting from time t , the program calls a cumulation sum function on the *ep_rewards* then subtract it from G (line 148). Thus *returns* would be storing the total discounted rewards for each time from time 0 to $T - 1$.

Then we will pass the list of states, actions, and the returns into the training step (line 154-157). Then tensorflow will use the state and the weights to generate a new μ and σ (line 63-93). Then it will compute the log probability of the actions given the generate μ and σ (line 102). The cost function used is $J(\theta) = -\sum [G_t \log \pi(A_t|S_t, \theta)]$. The program uses gradient descent to adjust the θ to minimize the cost function (line 105-107).

```

63 # 1 layer of hidden unit. Activation is ReLU
64 hidden = fully_connected(
65     inputs=x,
66     num_outputs=hidden_size,
67     activation_fn=tf.nn.relu,
68     weights_initializer=hw_init,
69     weights_regularizer=None,
70     biases_initializer=hb_init,
71     scope='hidden')
72
73 # use last layer of neural network as phi(a, s) (the feature)
74 # mu = phi(s, a)^T dot theta
75 mus = fully_connected(
76     inputs=hidden,
77     num_outputs=output_units,
78     activation_fn=tf.tanh,
79     weights_initializer=mw_init,
80     weights_regularizer=None,
81     biases_initializer=mb_init,
82     scope='mus')
83
84 # softplus is similar to ReLU. Activation function is g(x) = ln(1+e^x)
85 sigmas = tf.clip_by_value(fully_connected(
86     inputs=hidden,
87     num_outputs=output_units,
88     activation_fn=tf.nn.softplus,
89     weights_initializer=sw_init,
90     weights_regularizer=None,
91     biases_initializer=sb_init,
92     scope='sigmas'),
93     TINY, 5)
94
95 all_vars = tf.global_variables()
96
97 # use a Gaussian dist on pi because action is continuous
98 pi = tf.contrib.distributions.Normal(mus, sigmas, name='pi')
99 pi_sample = tf.tanh(pi.sample()), name='pi_sample')

```

Generating distribution on π

```

36 # if we have the w's and b's saved, load it. Otherwise initialize it
37 if args.load_model:
38     model = np.load(args.load_model)
39     hw_init = tf.constant_initializer(model['hidden/weights'])
40     hb_init = tf.constant_initializer(model['hidden/biases'])
41     mw_init = tf.constant_initializer(model['mus/weights'])
42     mb_init = tf.constant_initializer(model['mus/biases'])
43     sw_init = tf.constant_initializer(model['sigmas/weights'])
44     sb_init = tf.constant_initializer(model['sigmas/biases'])
45 else:
46     hw_init = weights_init
47     hb_init = relu_init
48     mw_init = weights_init
49     mb_init = relu_init
50     sw_init = weights_init
51     sb_init = relu_init

```

Weight (θ) initialization

```

115 track_returns = []
116 for ep in range(16384):
117     # reset the environment
118     obs = env.reset()
119
120     # generating all the states and actions and rewards
121     G = 0
122     ep_states = []
123     ep_actions = []
124     ep_rewards = [0]
125     done = False
126     t = 0
127     I = 1
128     while not done:
129         ep_states.append(obs)
130         env.render()
131         # pi_sample is the list of randomly generated probability
132         # then we use pi_sample to generate the list of actions
133         action = sess.run([pi_sample], feed_dict={x:[obs]})[0][0]
134         ep_actions.append(action)
135         obs, reward, done, info = env.step(action)
136         ep_rewards.append(reward * I)
137         G += reward * I # G is the total discounted reward
138         I *= gamma
139
140         t += 1
141         if t >= MAX_STEPS:
142             break
143     # done generating
144
145     if not args.load_model:
146         # G_t = total - culmulative up to time t
147         # set of all G_t's
148         returns = np.array([G - np.cumsum(ep_rewards[:-1])]).T
149         index = ep % MEMORY

```

States, actions, and rewards generation based on π_{sample}

Part 2

Question

Answer

Part 3

Question

Answer