

CSC411: Assignment # 4

Due on Tuesday, April 4, 2017

Ding, Si Chuang
1001165106

Shi, Yudian
999677862

Tuesday, April 4, 2017

1 Analysis of Bipedal Walker Algorithm

Note that REINFORCE uses the complete return from time t , which includes all future rewards up until the end of the episode. In this sense REINFORCE is a Monte Carlo algorithm and is well defined only for the episodic case with all updates made in retrospect after the episode is completed (like the Monte Carlo algorithms in Chapter 5). This is shown explicitly in the boxed pseudocode below.

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 θ
 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
 $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t|S_t, \theta)$

The vector $\frac{\nabla_{\theta} \pi(A_t|S_t, \theta)}{\pi(A_t|S_t, \theta)}$ in the REINFORCE update is the only place the policy parameterization appears in the algorithm. This vector has been given several names and notations in the literature; we will refer to it simply as the *eligibility vector*. The eligibility vector is often written in the compact form $\nabla_{\theta} \log \pi(A_t|S_t, \theta)$, using the identity $\nabla \log x = \frac{\nabla x}{x}$. This form is used in all the boxed pseudocode in this chapter. In earlier examples in this chapter we considered exponential softmax policies (13.2) with linear action preferences (13.3). For this parameterization, the eligibility vector is

$$\nabla_{\theta} \log \pi(a|s, \theta) = \phi(s, a) - \sum_b \pi(b|s, \theta) \phi(s, b). \quad (13.7)$$

As a stochastic gradient method, REINFORCE has good theoretical convergence properties. By construction, the expected update over an episode is in the same direction as the performance gradient.² This assures an improvement in expected performance for sufficiently small α , and convergence to a local optimum under standard stochastic approximation conditions for decreasing α . However, as a Monte Carlo method REINFORCE may be of high variance and thus slow to learn.

We applied this algorithm to the gridworlds in Examples 13.1 and/or 13.2 to obtain the results presented earlier or here in this section.

Figure 1: p.271 of Sutton & Barto

The policy function π_{θ} is implemented with a single hidden layer. The actions for the bipedal walker is continuous, so we have to use Gaussian distribution on π . Here, the hidden layer is passed to μ and σ .

If there is no weight θ saved previously, it will be initialized. For hidden, μ and σ layers, there are corresponding w's and b's. The initialization codes shown below keep the gradients scale roughly the same.

```
1 # check whether weights and biases are saved, load them. If not, initialize them.
2 if args.load_model:
3     model = np.load(args.load_model)
4     hw_init = tf.constant_initializer(model['hidden/weights'])
5     hb_init = tf.constant_initializer(model['hidden/biases'])
```

```

6     mw_init = tf.constant_initializer(model[ 'mus/weights '])
7     mb_init = tf.constant_initializer(model[ 'mus/biases '])
8     sw_init = tf.constant_initializer(model[ 'sigmas/weights '])
9     sb_init = tf.constant_initializer(model[ 'sigmas/biases '])
10  else :
11      hw_init = weights_init
12      hb_init = relu_init
13      mw_init = weights_init
14      mb_init = relu_init
15      sw_init = weights_init
16      sb_init = relu_init
17
18
19  # activation function: ReLU
20  # 1 layer of fully connected hidden unit
21  hidden = fully_connected(
22      inputs=x,
23      num_outputs=hidden_size ,
24      activation_fn=tf.nn.relu ,
25      weights_initializer=hw_init ,
26      weights_regularizer=None,
27      biases_initializer=hb_init ,
28      scope='hidden')

```

activation function :tanh

last layer of the neural network used as $\phi(a,s)(\text{feature})$ $\mu = \phi(s, a)^T$

```

1  mus = fully_connected(
2      inputs=hidden ,
3      num_outputs=output_units ,
4      activation_fn=tf.tanh ,
5      weights_initializer=mw_init ,
6      weights_regularizer=None,
7      biases_initializer=mb_init ,
8      scope='mus')

```

the sigma value is clipped when it is too small or too big

activation function: softplus $g(x) = \ln(1 + e^x)$

```

1  sigmas = tf.clip_by_value( fully_connected(
2      inputs=hidden ,
3      num_outputs=output_units ,
4      activation_fn=tf.nn.softplus ,
5      weights_initializer=sw_init ,
6      weights_regularizer=None,
7      biases_initializer=sb_init ,
8      scope='sigmas') ,
9      TINY, 5)

```

log probability of y given μ and σ . Lists of states, actions, and the returns will be passed into the training block (mention below in the for loop). Tensorflow will use the state and weights to generate new μ and σ and later compute the log probability of the actions based on the new μ and σ .

```

1  log_pi = pi.log_prob(y, name='log_pi')

```

returns 1 by (T-1) array for rewards (float). The cost function used here is $J(\theta) = -\sum[G_t \log \pi(A_t|S_t, \theta)]$. gradient decent is used to optimize the θ in order to minimize the cost.

```

1 Returns = tf.placeholder(tf.float32, name='Returns')
2 optimizer = tf.train.GradientDescentOptimizer(alpha)
3 train_op = optimizer.minimize(-1.0 * Returns * log_pi)

```

Here is where training data starts. In each iteration, the environment will be reset. Then, states, actions and rewards are generated from time 0 to time T. Randomly sample actions from the π normal distribution. Use `pi.sample()` to generate the corresponding action `pi_sample` and use it to generate new action, state and reward. Lists of `ep_actions`, `ep_states`, `ep_rewards` are used to store all corresponding information. The total discounted rewards is stored in G. Then we obtain G_t , the program calls a cumulative sum function on `ep_reward` and minus it from G. starting from time t. Finally, returns store the total discounted rewards for each iterations from time 0 to T-1.

```

1 for ep in range(16384):
2     # reset the environment
3     obs = env.reset()
4
5     G = 0
6     ep_states = []
7     ep_actions = []
8     ep_rewards = [0]
9     done = False
10    t = 0
11    I = 1
12    while not done:
13        ep_states.append(obs)
14        env.render()
15        # pi_sample() is a list store randomly generated probabilities
16        # assign pi_sample to action
17        action = sess.run([pi_sample], feed_dict={x:[obs]})[0][0]
18        ep_actions.append(action)
19        obs, reward, done, info = env.step(action)
20        ep_rewards.append(reward * I)
21        # G is total discounted reward
22        G += reward * I
23        I *= gamma
24        t += 1
25        if t >= MAX.STEPS:
26            break
27    if not args.load_model:
28        # returns = G_t = total discounted reward - cumulative up until time t
29        returns = np.array([G - np.cumsum(ep_rewards[:-1])]).T
30        index = ep % MEMORY
31        - = sess.run([train_op],
32                    # ep_states includes all states from S_0 to S_{T-1}
33                    # ep_actions includes all actions from A_0 to A_{T-1}
34                    feed_dict={x:np.array(ep_states),
35                               y:np.array(ep_actions),
36                               # returns include all returns(G_t) from t = 0 until t
37                    = T
38                    Returns:returns })

```

2 Part 2, Changes To Code

In part 1, we used multiple layers: hidden layer which computes the feature vector x . It is then been passed to the other layer which computes μ and the other computes σ . Based on the outputs we get to generate π and further use π found to compute action.

In part 2, it is the simplified version based on the method used in part 1. Instead of using multiple layers, we only used one layer here. We first create the environment 'CartPole-v0'. The activation function used is "sigmoid". An output array of size 2 is generated, of 0 and 1, sampling from a Bernoulli distribution with the probability determined by the observation and weights. Each of the two elements from the policy output represent whether the cart should go left or right. The cart will receive a input force to the left if the network sampled a 0 for moving right and a 1 for moving left. Otherwise the cart will receive a pump to the right. This is implemented as a argmax on the output vector.

2.1 Detailed Changes to Policy Related Code

This section lists the relevant changes to the example we required to make the policy network function. Some other less relevant code, such as plotting and book keeping, that are does not influence funtion of the policy network are omitted.

Changed the learning rate to 0.001 and gamma to 0.99, as per instructed in the assignment. We also found that these values did produce much better result than others.

```
1 # Parameters for single layer NN as policy network
2 hidden_size = 8
3 alpha = 0.001
4 TINY = 1e-8
5 gamma = 0.99
```

Changed input and output nodes to the correct dimension

```
1 # Tensorflow input and output initialization
2 # Input features: 4, output features 2
3 x = tf.placeholder(tf.float32, shape=(None, 4), name='x')
4 y = tf.placeholder(tf.float32, shape=(None, 2), name='y')
```

Modified the nerual net structure to only a single softmax layer, with sigmoid activation function. The assignemnt suggested that we use this structure for the policy network. The sigmoid activation function layer is equivlent to softmax.

```
1 # Only one lay network is required, softmax
2 sigmas = tf.clip_by_value(fully_connected(
3     inputs=x,
4     num_outputs=2,
5     activation_fn=tf.nn.sigmoid,
6     weights_initializer=sw_init,
7     weights_regularizer=None,
8     biases_initializer=sb_init,
9     scope='sigmas'),
10     TINY, 5)
```

Modified the probability distribution for pi. Since the action here is descrete the bernulie distribution is more appropiate space to sample from. Sampling from this distribution will give us a

element output, representing if the network decide if we should push left, right, or undecided. In this case of this implimentation, the cart force act to the right by default if the sampled outputs are both 1 or 0.

```
1 # Layers of network
2 # Use bernulie distribution here instead since actions are descrete
3 pi = tf.contrib.distributions.Bernoulli(p = sigmas, name='pi')
4 pi_sample = pi.sample()
5 log_pi = pi.log_prob(y, name='log_pi')
```

The while loop running the enviornment is modified to handle the previous changes made to the policy network.

```
1 while not done:
2     ep_states.append(obs)
3     #env.render()
4     action = sess.run([pi_sample], feed_dict={x:[obs]})[0][0]
5
6     obs, reward, done, info = env.step(np.argmax(list(action)))
7
8     ep_actions.append(action)
9     ep_rewards.append(reward * I)
10    G += reward * I
11    I *= gamma
12
13    t += 1
14    if t >= MAX_STEPS:
15        break
```

3 Training Cartpole

3.1 part 3a

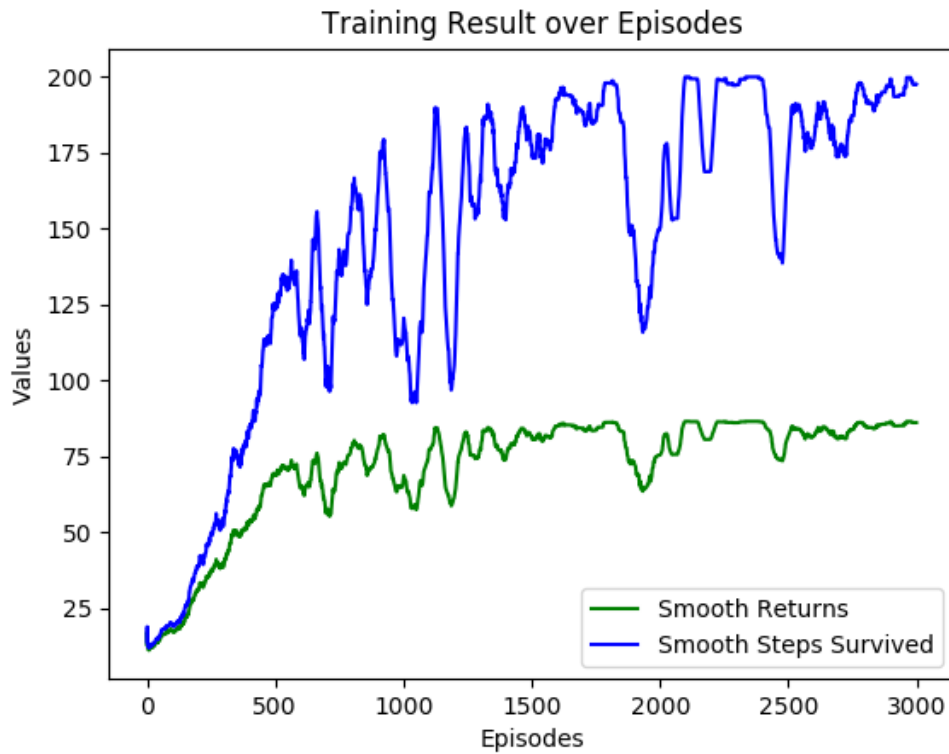
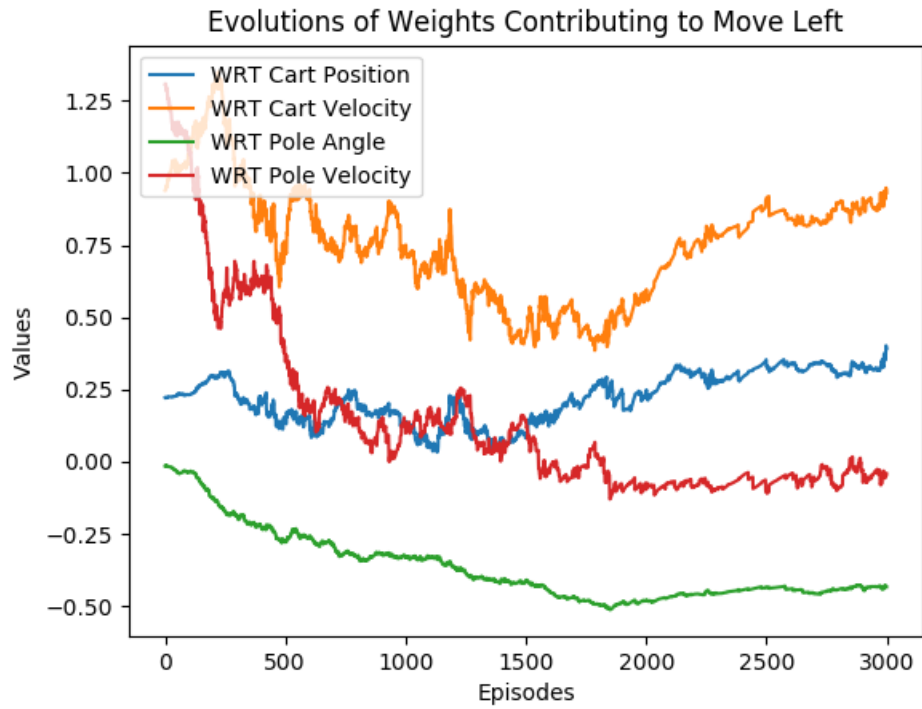
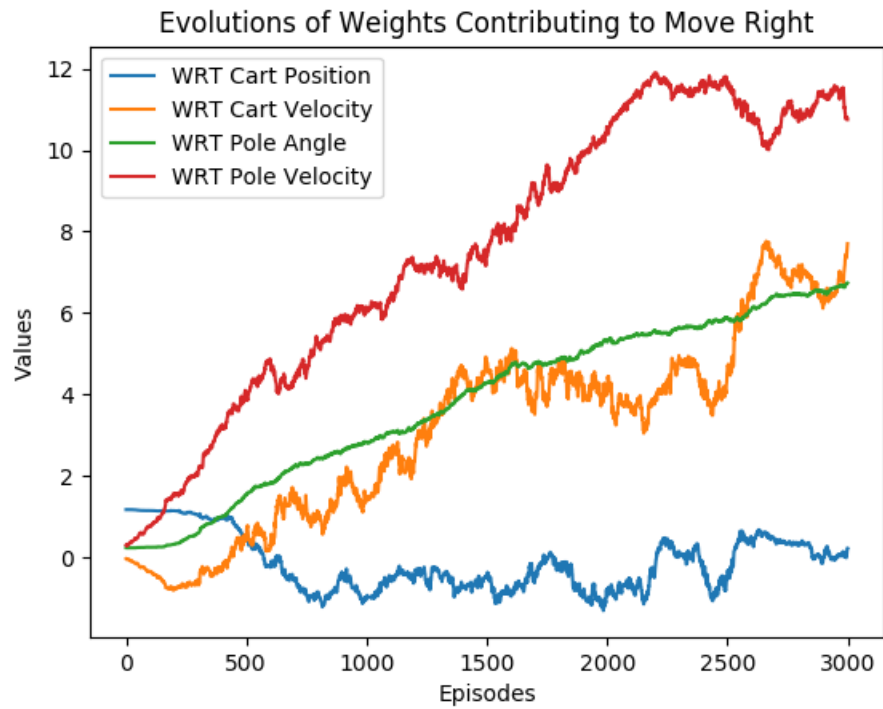


Figure 2: Steps Survived and Returns Over Training Episodes



(a) Weights for Moving Left



(b) Weights for Moving Right

Figure 3: Evolution of Weights Over Episodes

3.2 Part 3b, Why Weights Work

By the nature of how cartpole is implemented, there is no “No Action” move. i.e. the cart have to neither move left or move right. Also, from a controls perspective, we know that since we only have the position and velocity information from cartpole, we would expect some kind of PD controller to be implemented by the policy network. Looking at the weights, it makes a huge amount of sense:

The weights for moving left is very small compared to the ones for moving right. This makes a lot of sense since the cart can only do one of two actions. If it is not moving right then it will automatically be moving left. So only one of the output nodes need to have large weights that are trained to the system and the other one can essentially be passive. In this case, the decision ended up being neither one of “move right” or “not move right”.

Focusing on the weights toward direction the cart to move right, we noticed the properties that we expect out of a PD controller.

- The weight is largest with the pole velocity, the cart is most likely to move right if it system sees the pole is falling over toward the right. Weights for the pole velocity if about double compared to cart velocity and pole position. This makes sense because the cart need to move twice as much to influence the pole velocity as the center mass of the pole is half way up, accelerating the cart by a unit velocity will accelerate the pole by half that much.
- Also, the weights are positively correlated to pole position, so if the pole is leaning toward the right side of the cart, the cart must move to the right.
- Weights are also positively correlated to the care velocity, since if the cart is moving to the right with some velocity, the cart need to be accelerated more to correct a pole leading right.
- Finally, the most telling and intuitive weight is the weight corresponding to the cart position. Its obvious from the setup of the environment, that the absolute position of the cart have no effect on balancing the pole, as long as the cart do not run off the edge of its rail. We can clearly see that the policy network was able to infer this fact, and adjusted the weights corresponding to the car's position to essentially be zero.