# CSC411 Project 4: Reinforcement Learning with Policy Gradients

Dawn Duan and Shamama Khattak

April 6, 2017

## 1  Part 1

The pseudocode on p. 271 of Sutton Barto is as follows:

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

The precise explanation of why the code corresponds to the pseudocode is as follows.

1. **Declare all the input variables, i.e. the policy variable $\pi$, action variable $a$, state variable $s$.**

   $\pi_\theta(a|s)$ takes state s at time t, or at the current state, as input features and outputs the action the agent should pursue at the next state, $s_{t+1}$

   ```
   action = sess.run([pi sample], feed dict={x:[obs]})[0][0]
   ```

   Here, a continuous Gaussian policy is adopted here because the metion of the bipedal walker is continuous.

   ```
   pi = tf.contrib.distributions.Normal(mus, sigmas, name='pi')
   ```

2. **Initialization policy weights $\theta$.**

   ```
   1   weights_init = xavier_initializer(uniform=False)
   2
   3   if args.load_model:
   4       model = np.load(args.load_model)
   5       hw_init = tf.constant_initializer(model['hidden/weights'])
   6       hb_init = tf.constant_initializer(model['hidden/biases'])
   7       mw_init = tf.constant_initializer(model['mus/weights'])
   8       mb_init = tf.constant_initializer(model['mus/biases'])
   9       sw_init = tf.constant_initializer(model['sigmas/weights'])
   10      sb_init = tf.constant_initializer(model['sigmas/biases'])
   11  else:
   12      hw_init = weights_init
   13      hb_init = relu_init
   14      mw_init = weights_init
   15      mb_init = relu_init
   16      sw_init = weights_init
   17      sb_init = relu_init
   18
   ```

```
19    try:
20        output_units = env.action_space.shape[0]
21    except AttributeError:
22        output_units = env.action_space.n
23
24    hidden = fully_connected(
25        inputs=x,
26        num_outputs=hidden_size,
27        activation_fn=tf.nn.relu,
28        weights_initializer=hw_init,
29        weights_regularizer=None,
30        biases_initializer=hb_init,
31        scope='hidden')
32
33    mus = fully_connected(
34        inputs=hidden,
35        num_outputs=output_units,
36        activation_fn=tf.tanh,
37        weights_initializer=mw_init,
38        weights_regularizer=None,
39        biases_initializer=mb_init,
40        scope='mus')
41
42    sigmas = tf.clip_by_value(fully_connected(
43        inputs=hidden,
44        num_outputs=output_units,
45        activation_fn=tf.nn.softplus,
46        weights_initializer=sw_init,
47        weights_regularizer=None,
48        biases_initializer=sb_init,
49        scope='sigmas'),
50        TINY, 5)
```

3. **Start the loop.**

```
1    for ep in range(16384):
```

The range, 16384 is the large arbitrary value which we want to consider as the last episode.

4. **Generate an Episode** $S_0, A_0, R_1, ..., S_{T-1}, A_{T-1}, R_T$, **following** $\pi(.|.,\theta)$**.** The declaration of the terms is presented below, where the subscript means the time-step.

```
1    # initilization
2        obs = env.reset()
3
4        G = 0
5        ep_states = []
6        ep_actions = []
7        ep_rewards = [0]
8        done = False
9        t = 0
10       I = 1
```

5. **For each step of the episode** $t = 0, ..., T-1$**.**

The declaration of the variable, time-step $t$ is shown below. The last time-step, i.ei. $T$ in the pseudocode, corresponds to the MAX STEPS defined from the *env*. That is, when the done flag is set to True, whichever comes first:

```
1    MEMORY=25
2    MAX_STEPS = env.spec.tags.get('wrapper_config.TimeLimit.max_episode_steps')
```

```
3
4    t = 0
5    while not done:
6        ...
7        t += 1
8        if t >= MAX_STEPS:
9            break
```

6. **Get $G_t$, is the actual total reward in each episode starting from time t, i.e. discounted.**

$$G_t = \sum_{i=t} \gamma^i r_i$$

In this case, starting from time t $= 0$

```
1    gamma = 0.98
2    '''iterate through all the episode'''
3    G = 0
4    t = 0
5    I = 1 # discounted rate
6
7    obs, reward, done, info = nev.step(action) # gives out the reward value at each action step, tell you weather goal stae
8    G += reward * I #update the actual reward at each step
9    I *= gamma #update discounted rate
10   t += 1 #update time step
```

Here a difference between the pseudo code and the real implementation is presented. That is, the actual reward G is supposed to get updated at each time-step for each episode. In the actual implementation, the actual reward is only override in each episode. For the actual reward in each time-step, it is implemented using np.cumsum() to compute the reward $ep_rewards$, for time-step 0 up to time-step T.

```
1        # Updates for the rewards
2        if not args.load_model:
3            returns = np.array([G - np.cumsum(ep_rewards[:-1])]).T
```

7. **Update $\theta$ using gradient ascend, i.e. the following formula:**

$$\theta \leftarrow \theta + \alpha\gamma^t G_t \nabla_\theta log(\pi_\theta(A_t|S_t))$$

The learning rate $\alpha$ is specified as 0.01 here.

```
1    alpha = 0.01
2
3    optimizer = tf.train.GradientDescentOptimizer(alpha)
4    train_op = optimizer.minimize(-1.0 * Returns * log_pi) # the update to theta
5    # -1.0 because it's gonna be adding for gradient descend
6
7    # the relevant values are feed into the tensor per episode
8        if not args.load_model:
9            returns = np.array([G - np.cumsum(ep_rewards[:-1])]).T
10           index = ep % MEMORY
11
12
13           _ = sess.run([train_op],
14                   feed_dict={x:np.array(ep_states),
15                              y:np.array(ep_actions),
16                              Returns:returns })
```

The gradient is computed from Policy Gradient Theorem. It consists three part:

(1) $G_t$ is computed from the previous step, Step5

(2) $\nabla_\theta log(\pi_\theta(A_t|S_t))$ is presented as follows:

```
1  pi = tf.contrib.distributions.Normal(mus, sigmas, name='pi')
2  pi_sample = tf.tanh(pi.sample(), name='pi_sample')
3  log_pi = pi.log_prob(y, name='log_pi')#insert tensor y
4
5  train_op = optimizer.minimize(-1.0 * Returns * log_pi) # the update to theta
6
7  # loop through eps
8  action = sess.run([pi_sample], feed_dict={x:[obs]})[0][0]
9  obs, reward, done, info = env.step(action) #obs are the states
```

# 2  Part 2

In this section, we are required to write a reinforcement learning algorithm for the "CartPole-v0" environment. In order to accomplish this task, our group made the modifications detailed below to the REINFORCE reinforcement learning algorithm used for the "BipedalWalker-v2" environment.

Previously, in the "BipedalWalker-v2" env, we computed the policy function, $\pi_\theta$, as a single hidden layer neural network. Since the "CartPole-v0" environment defines two discrete actions (move left or move right), to obtain a good policy function and make REINFORCE work, we had to do the following:

1. **Change from Normal to Bernoulli Distribution**

   Since we are using reinforcement learning, we want to use a stochastic policy that, given an action $a$ and a state $s$, gives us the probability of taking action $a$ from state $s$ at time step $t$. That is we aim to have a policy, $\pi_(a|s) = P(A_t = a, S_t = s)$

   Previously, since we had a continuous set of actions, we used the gaussian policy. Now that we have a discrete set of actions, we must now use a softmax probability.

   More formally, we aim to change from: $a \sim N(f_\theta(s))$ to $a \sim Bernoulli(f_\theta(s))$. That is, we need to change from a Normal distribution to a bernoulli distribution.

   In code this change corresponds to (and resulted in the change of):

```
1  # Define policy function to be a Bernoulli distribution
2  pi = tf.contrib.distributions.Bernoulli(p=dir_probabilities, name="pi")
3
4  pi_sample = pi.sample()
```

2. **Change the function $f_\theta(s)$**

   We also changed the function $f_\theta(s)$ so that instead of being a neural network with one hidden layer, it is now a fully connected network with no hidden layers. It uses the softmax activation function and has 2 output units which correspond to the actions that we can take. The outputs represent the probability of taking action left or taking action right. Our network takes in x which is a NONE (variable) X 4 matrix. We chose to use a NONE X 4 matrix because we have 4 numbers in our observation space as found by executing the command: env.observation_space.shape[0]. The code corresponding to the change is shown below:

```
1  # Gets observation space set size
2  Y_UNITS = 2
3  input_shape = env.observation_space.shape[0]
```

```
4   NUM_INPUT_FEATURES = 4
5   x = tf.placeholder(tf.float32, shape=(None, NUM_INPUT_FEATURES), name='x')
6   y = tf.placeholder(tf.float32, shape=(None, Y_UNITS), name='y')
7
8   # Obtain probabilities of taking action left and right
9   dir_probabilities = fully_connected(
10      inputs=x,
11      num_outputs=output_units,
12      activation_fn=tf.nn.softmax,
13      weights_initializer=w_init,
14      weights_regularizer=None,
15      biases_initializer=b_init,
16      scope='dir_probabilities')
```

3. **Other changed lines:**

The following lines were changes so that the training part would work with the changes we made above:

```
1   1. MAX_STEPS = 200
2   2. obs, reward, done, info = env.step(action[0])
3   3. _ = sess.run([train_op],
4                      feed_dict={x:np.array(ep_states),
5                                 y:np.vstack(np.array(ep_actions)),
6                                 Returns:returns })
```
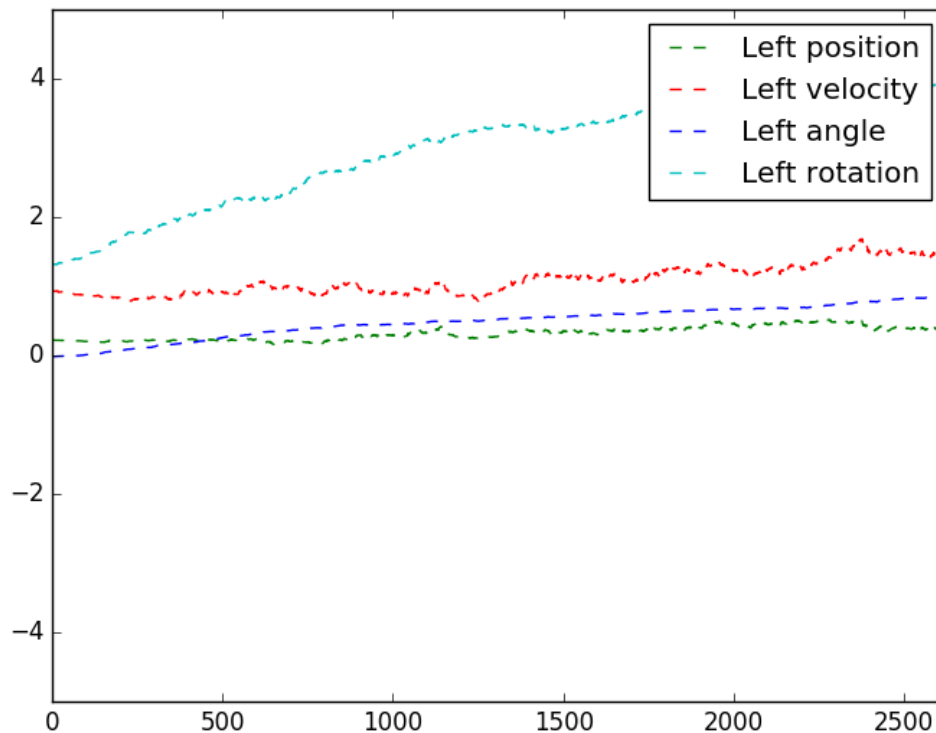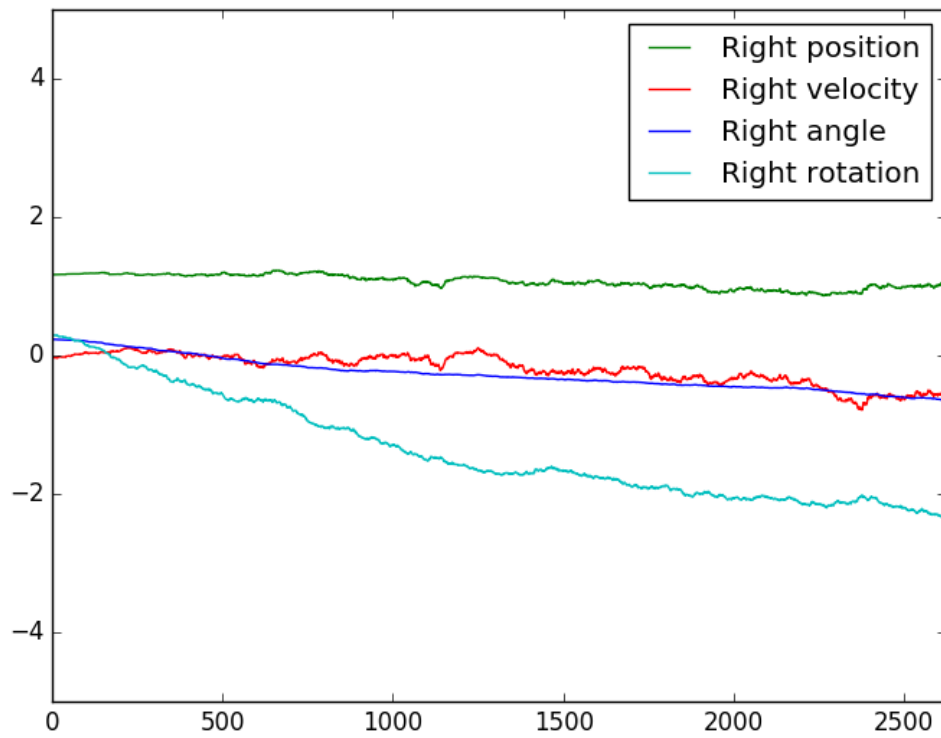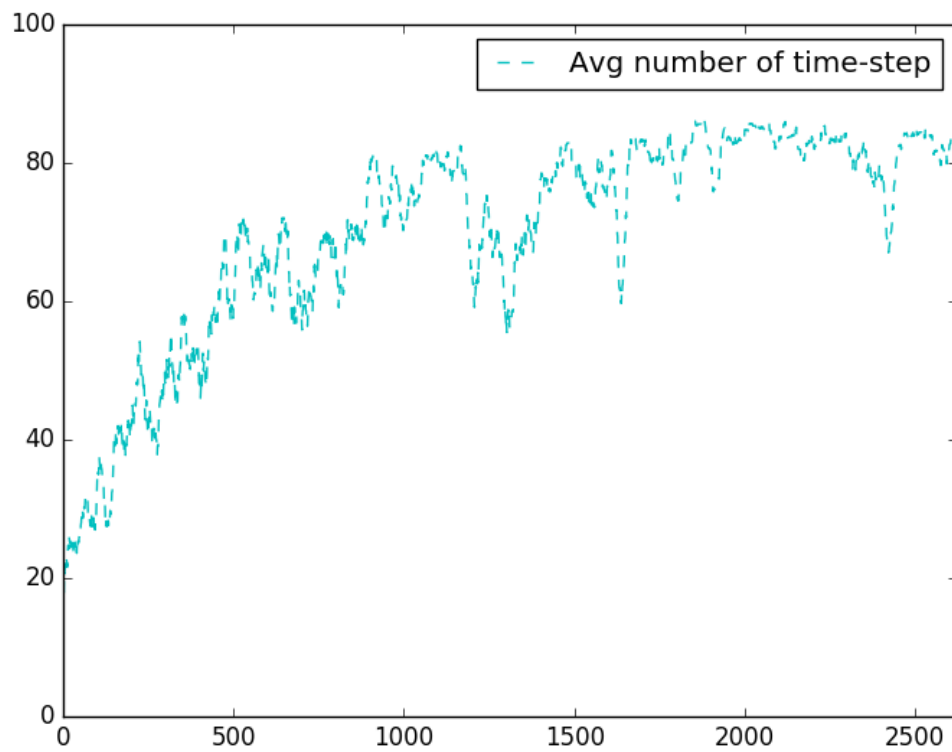
# 3 Part 3

## 3.1 Print-out of weights

The printout that shows how the weights of the policy function changed for both output units are as follows:

The printout that shows how the average number of time-steps per eipsode (over e.g. the 25 last episodes) changed as we trained the agent is as follows:

## 3.2 Explanation of weights

The four input dimensions correspond to:

1. The position of the cart

2. The velocity of the cart

3. The angle of the pole

4. The rotation rate of the pole

The final weights that we obtain are:

$$[[ \ 0.20637971, \quad 1.18371701],$$
$$[ \ 0.95413405, \quad -0.05625207],$$
$$[ \ 0.37281665, \quad -0.15437123],$$
$$[ \ 2.46482682, \quad -0.85481673]]$$

In order to understand why the weights listed above make sense, we must first consider what each component of the 2 x 4 matrix corresponds to. Each of the 2 columns represents the weights for each of the output units. Each row corresponds to a feature of our observation state space. Then, using this line of reasoning, it can be inferred that the weights at index 1, 1 is the weight that corresponds to the position of the cart for the first output unit.

Let's consider each of the features of the observation space individually in order to deduce why these values make sense:

1. **Feature: Position of the cart**

   The weights for the position of the cart is: $[0.20637971 \quad 1.18371701]$

   Consider that the position on the cart is on a scale from -5 to 5, where the negative axis indicates that the cart is on the left of the center and positive implies that we are to the right of the center. Then the position weight **multiplied** by the position state for output 1 will be higher than the position state times the position weight for output 2 when the cart is to the left of the center. When the weight is to the left of the center, the position described by the input is negative. The output 1 weight is smaller (0.20637971) than the output 2 weight (1.18371701) and the effect of multiplying a large negative number by a small positive number results in a smaller outcome than multiplying a large negative number by a large positive number.

   Thus, we can say that output unit 1 is favoured (i.e., larger than output unit 2) when the cart is to the left and output unit 2 is favoured when the cart is to the right.

2. **Feature: Velocity of the cart**

   The weights for the velocity of the cart is: $[ \ 0.95413405, -0.05625207]$

   These weights imply that the velocity of the cart will be higher in output unit 1 when the input state velocity is positive (since the weight is positive and a positive number multiplied by a positive number is greater than a positive number multiplied by a negative number). Conversely, output unit 2 will be favoured when the input state velocity of the cart is negative (implying that the cart is moving towards the left).

3. **Feature: Angle of the pole**

   The weights for the angle of rotation is: $[ \ 0.37281665, -0.15437123]$

   When the input state angle of the pole is positive (tilting to the right) we will favour output unit 1 (by similar logic to the weights above). When the the input state angle of the pole is negative (i.e., tilting to the left), we will favour output unit 2.

4. **Feature: Rotation rate of the pole**

   The weights for the angle of rotation is: $[2.46482682, -0.85481673]$

When the input state's rotation rate of the pole is negative (i.e., the pole is falling to the left) that means we are likely to favour output unit 2. When the input state's rotation rate of the pole is positive (falling to the right), we are very strongly likely to favour output unit 1.

Putting all of this together, what we find is that output unit 1 is more likely to be favoured (greater than output unit 2) when the input for our weight is as follows: the cart is to the left of the center, moving with a positive velocity to the right with the pole leaning to the right, and falling to the right. Intuitively, the combination of these parameters as weights for output unit 1 makes sense. In the steady state, we want to occasionally try to oppose the action given so that we can stabilize ourselves. Using similar logic, we can come up with a similar deduction for output unit 2.