

# **CSC411: Project 4**

Due on Friday, March 31, 2017

**Katie Datsenko, Loora Zhuoran Li**

April 3, 2017

## Part 1

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

Input: a differentiable policy parameterization  $\pi(a|s, \theta), \forall a \in A, s \in S, \theta \in R^n$

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  $t$

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

We would like to show the code in `bipedal-reinforce.py` corresponds with the pseudocode above.

#### 1. Input and initialization, policy learning

A policy function  $\pi$  takes in a current state  $\mathbf{s}$ , and outputs the action  $\mathbf{a}$  the agent should take. For stochastic policy,  $\pi(a|s) = P(A_t = a | S_t = s)$

In `bipedal-reinforce.py`, a single-hidden-layer neural network is used to implement the policy function.

$\mathbf{x}$  – Input for this neural network

    a state vector of dimension 24 (24 features for one state  $\mathbf{s}$ )

$\mathbf{y}$  – Output for this neural network

    probability for taking each possible action given current state  $\mathbf{s}$ , number of actions = number of output units

```
input_shape = env.observation_space.shape[0]
NUM_INPUT_FEATURES = 24 #the state vector
x = tf.placeholder(tf.float32, shape=(None, NUM_INPUT_FEATURES), name='x')
y = tf.placeholder(tf.float32, shape=(None, output_units), name='y')
```

Referring to the first two lines of the pseudocode.

Before implementing the network structure, the policy weights and bias terms are initialized randomly.

```
weights_init = xavier_initializer(uniform=False)
relu_init = tf.constant_initializer(0.1)

if args.load_model:
5   model = np.load(args.load_model)
   hw_init = tf.constant_initializer(model['hidden/weights'])
   hb_init = tf.constant_initializer(model['hidden/biases'])
   mw_init = tf.constant_initializer(model['mus/weights'])
   mb_init = tf.constant_initializer(model['mus/biases'])
10  sw_init = tf.constant_initializer(model['sigmas/weights'])
   sb_init = tf.constant_initializer(model['sigmas/biases'])
else:
   hw_init = weights_init
   hb_init = relu_init
15  mw_init = weights_init
   mb_init = relu_init
   sw_init = weights_init
```

```

        sb_init = relu_init
    try:
20     output_units = env.action_space.shape[0]
    except AttributeError:
        output_units = env.action_space.n

```

In terms of structure for the network,

— "hidden" is a fully connected layer, with observation features for states as inputs, hidden units as outputs, relu function as the activation function, weights initialized using xavier initializer/given constant, bias terms initialized using relu initializer/given constant.

In this case, the actions are continuous, thus it is reasonable to use Gaussian distribution for the actions, centered around  $\phi(s, a)^T \theta$ , which is a  $\sim N(\phi(s, a)^T \theta, \sigma^2)$

— "mus" is a fully connected layer, with hidden units as inputs,  $\phi(s, a)^T \theta$  as outputs, tanh function as the activation function, weights initialized using xavier initializer/given constant, bias terms initialized using relu initializer/given constant.

— "sigma" is a fully connected layer, also with hidden units as inputs, standard deviation  $\sigma$  as outputs, softplus function as the activation function, weights initialized using xavier initializer/given constant, bias terms initialized using relu initializer/given constant.

```

hidden = fully_connected(
    inputs=x, #inputs x
    num_outputs=hidden_size,
    activation_fn=tf.nn.relu,
5    weights_initializer=hw_init,
    weights_regularizer=None,
    biases_initializer=hb_init,
    scope='hidden')

10 #Gaussian distribution policy for the actions
# mus = phi(s, a)^T times theta
mus = fully_connected(
    inputs=hidden, #hidden layer
    num_outputs=output_units,
15    activation_fn=tf.tanh,
    weights_initializer=mw_init,
    weights_regularizer=None,
    biases_initializer=mb_init,
    scope='mus')

20 #the sigmas of the gaussian distribution,
sigmas = tf.clip_by_value(fully_connected(
    inputs=hidden, #units of hidden layer
    num_outputs=output_units,
25    activation_fn=tf.nn.softplus,
    weights_initializer=sw_init,
    weights_regularizer=None,
    biases_initializer=sb_init,
    scope='sigmas'),
30    TINY, 5)

```

$$a \sim N(\phi(s, a)^T \theta, \sigma^2)$$

For a certain state, probabilities  $\pi$  are computed for continuous actions based on the state input  $x$  with 24 features.  $\pi$ -sample is a randomly generated action based on probabilities in  $\pi$ , this is called later in the loop.

```
all_vars = tf.global_variables()
pi = tf.contrib.distributions.Normal(mus, sigmas, name='pi')
pi_sample = tf.tanh(pi.sample(), name='pi_sample')
```

## 2. Repeat Forever: loop for infinite number of episodes

Reset the environment for a new episode, following the policy with current  $\theta$ , generate all the states and actions  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ .

$G$  is the total (discounted) reward, starting from time  $t=0$

$t$  is the timestep

$I$  is the coefficient for reward at timestep  $t$  (this is a discounted reward system)

```

# reset the environment
obs = env.reset() #root state
# following the policy with current theta
# generate all the states and actions and rewards
5  G = 0
   ep_states = []
   ep_actions = []
   ep_rewards = [0]
   done = False
10  t = 0
   I = 1
```

Now we would like to analyze how  $G_t$  and  $\theta$  are updated

If this episode has not finished yet, add observation vectors( $x$  inputs) to the state list. Based on state observations from input  $x$ , pick a single random action at time step  $t$ . Retrieve information on observations, reward, whether episode finished or not and info for taking this action. Add discounted reward for taking this action  $a$  to the reward list. **Update total discounted reward  $G$  by adding the discounted reward at the current time step to its original value,  $G = G + reward_t$ .** Update the reward discount coefficient by multiplying  $\gamma$  to its original value. Update time step by adding 1. Repeat above steps until the average number of time steps per episode exceeds the constant value MAX STEP.

```

while not done:
    ep_states.append(obs)
    env.render()
    # pick a single random action for time step t, based on state obs -- input x
    # pi_sample is a randomly generated action based on probabilities in pi
5    action = sess.run([pi_sample], feed_dict={x:[obs]})[0][0]
    ep_actions.append(action)
    obs, reward, done, info = env.step(action)
    ep_rewards.append(reward * I) #R_t - reward for taking action a
10    # G is the total discounted reward, starting from time t=0
    #R0 + gamma*R1 + gamma^2*R2 + ... + gamma^(MAX_STEPS - 1)*R(N-1)
```

```

15      G += reward * I
      I *= gamma #rewards farther away in time count less (greedy approach)
      #gamma between 0 and 1
      t += 1
      if t >= MAX_STEPS:
          break

```

For each step of the episode, if not done,  $\theta$  is also updated. This corresponds to the last line in the pseudocode,  $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla \theta \log \pi(A_t | S_t, \theta)$

The codes in the next box show details of how the stochastic gradient is updated.

Take the log probability of the output  $\pi$  for each state in this current episode,  $y$  is the output of neural network ie. the probability for taking each possible action given current state  $s$ .

Initialize Returns in order to record all  $G_t$  from  $t = 1$  to  $T$

Multiply each reward value  $G_t$  with the log of probability of the action  $A_t$  taken at state  $S_t$  for time step  $t$ , then take summation over all values for each time step  $t$  within an episode.

The goal is to maximize expected returns ( $\text{Returns} \times \log(\pi)$ ), this represents rewards from policy with current  $\theta$ . Comparing with the pseudocode, instead of doing an update for every time step within an episode, it computes the cumulative gradients for all time steps together (the gradient of a sum is equal to the sum of the gradients of the individual terms). Then, multiply the expected returns by -1, thus we need to find the minimum of this function ( $-1 \times \text{Returns} \times \log(\pi)$ ). Apply gradient descent technique with Tensorflow to compute the gradient corresponding to the  $\theta$  for updates.

```

# log of prob of actions from A_0 to A_T-1
log_pi = pi.log_prob(y, name='log_pi') #log prob(p16)

#Returns contains all the G_t's from t=1 to t=T
5 Returns = tf.placeholder(tf.float32, name='Returns')
optimizer = tf.train.GradientDescentOptimizer(alpha)
train_op = optimizer.minimize(-1.0 * Returns * log_pi)

sess = tf.Session()
10 sess.run(tf.global_variables_initializer())

```

If the model not loaded: for each index  $i$  of ep rewards, take a cumulative sum of entry at index  $i$  with all the entries before it.  $G_t$  is the total discounted reward starting from time step  $t$ .  $G_t = \text{total rewards value} - \text{cumulative sum up towards time step } t$ . Returns is a set of all  $G_t$  from  $t = 1$  to  $T$  Run tensorflow Session, use gradient descent to maximize expected rewards thus update  $\theta$ ; record all the states from  $s_0$  to  $s_{T-1}$ ; record all the actions from  $A_0$  to  $A_{T-1}$ ; record all the discounted returns  $G_t$  from  $t = 1$  to  $T$ .

```

    if not args.load_model:
        returns = np.array([G - np.cumsum(ep_rewards[:-1])]).T
        index = ep % MEMORY

5    _ = sess.run([train_op],
                  feed_dict={x:np.array(ep_states),
                              y:np.array(ep_actions),
                              Returns:returns })

```

## Part 2

### 1. Environment Setup

We added our code right after the class *CartPoleEnv*, and set up the environment with these two lines:

```
env = gym.make('CartPole-v0')
```

We added the line

```
np.random.seed(0)
```

as part of our random seed settings because we use `np.random` further in the code (when generating an action via bernoulli sample).

### 2. Policy Function

We removed the previous policy function model which was a single-hidden-layer network and replaced it with a single fully-connected softmax layer.

We have variables  $x$ ,  $params$  with dimensions  $(T \times 4)$  and  $(4 \times 2)$  respectively. Variable  $x$  represents the list of States  $S$  of the cartpole, for a single episode. Each individual state is encoded by 4 features (thus, `NUM_INPUT_FEATURES` is updated to 4), relevant to the internal state representation as dictated by the *CartPoleEnv* class. Variable  $params$  represents our theta parameters for the policy function, they are the weights of the fully connected single layer. In the Cart Pole task, only two actions are possible, applying a force pushing left, and applying a force pushing right. Thus our policy function, which takes a state as input and outputs the probability of an action, has two outputs: one for action left and one for action right. Our corresponding single layer network has two output units, which are passed through a softmax activation function, which will eventually learn to compute the "probabilities" of an action, through our policy reinforcement training algorithm. This justifies the  $(4 \times 2)$  dimension of our theta parameter matrix  $params$ .

```
#the number of features of the state vector
NUM_INPUT_FEATURES = env.observation_space.shape[0]

params = tf.get_variable("thetas", shape=[NUM_INPUT_FEATURES, output_units])
5 bias = tf.get_variable("bias", shape=[output_units])

#the state vector input
x = tf.placeholder(tf.float32, shape=(None, NUM_INPUT_FEATURES), name='x')
#y - At (actions selected from A_0 to A_T-1 via bernoulli)
10 y = tf.placeholder(tf.int32, shape=(None), name='y')

linear_layer = tf.matmul(x, params)+bias
pi = tf.nn.softmax(linear_layer)
```

### 3. Stochastic gradient update

We also compute each stochastic estimation of the total expected reward based on an episode of the policy in the following way: we take the log of the softmax probabilities computed for every state in the episode (stored in Variable  $pi$ ). Variable  $y$  stores the list of actions (one action per time step in the episode), 0 is assigned as the left action, 1 as the right. We keep the probabilities of only the actions

that actually occurred in the generated episode by converting  $y$  to a one-hot representation and doing elementwise multiplication (Variable  $act\_pi$  with modifications to the line from the handout). Then we multiply each reward value  $G_t$  (stored in Returns) with the log of the probability of the action  $A_t$  taken at the state  $S_t$  for time step  $t$  and sum over all values for every time step  $t$  of an episode. This corresponds to the pseudocode on the slides, with the exception that we are trying to do the updates for an episode all at once, versus for each time step individually in an iteration of the loop and combining updates as we iterate (the gradient of a sum is equal to the sum of the gradients of the individual terms). To turn our Reward Expectation Maximization problem into a minimization problem, we multiply the expected episode reward by  $-1$ . The Tensorflow gradient descent algorithm will take care of the computation of the gradient with respect to  $\theta$  for this update.

```

# log of prob of actions A_0 and A_1 for every state for episode
log_pi = tf.log(pi) #2 prob values for every state in vector of time steps

5 #We make a one-hot vector of 0/1 actions stored in y
  #with a one at the action we want to increase the probability of.

#1. log_pi - should be (T x 2), 2 action probabilities for every time step
#2. tf.one_hot(y, 2, axis=1) - every row corresponds to one time step T,
10 #columns are left & right probabilities - (T x 2)

#tf.multiply - Returns x * y elementwise, will set the prob of action != At to 0
#tf.reduce_sum - squeezes vector values across two rows into flat vector,
#act_pi is probability of action A_t for every time step t
15 act_pi = tf.reduce_sum(tf.multiply(log_pi, tf.one_hot(y, 2, axis=1)),
  reduction_indices=[1])

#Returns contains all the G_t's from t=1 to t=T
20 Returns = tf.placeholder(tf.float32, name='Returns')
optimizer = tf.train.GradientDescentOptimizer(alpha)
train_op = optimizer.minimize(-1.0 * Returns * act_pi)

vect = act_pi * Returns #elementwise multiplication
25 #for gradient update for entire episode we take the sum
  #Multiply by -1 because we want to maximize JavV reinforcement reward by following policy theta
loss = -tf.reduce_sum(vect)
train_op = optimizer.minimize(loss)

```

#### 4. Training Parameters

We changed the values of  $\alpha$  and  $\gamma$  to be 0.0001 and 0.99 respectively, and compute the average number of time-steps per episode, stopping the reinforcement algorithm when the average is at 50 or higher. MAX\_STEPS (max limit of time steps possible in an episode) is given by the environment constant which is set at 200 by default.

```

alpha = 1e-6
gamma = 0.99
MAX_STEPS = env.spec.tags.get('wrapper_config.TimeLimit.max_episode_steps')

```

#### 5. Sampling action $A_t$

The last important change to the code is how we sample actions for each time step of an episode from our policy function with the current theta. For Bipedal reinforcement, the actions were continuous values; thus the probabilities of actions were modelled using a normal distribution, and a sample action for a state came from a Normal Distribution with a particular  $\mu$  and  $\sigma$  (computed using a network trained on state input). In the Cartpole problem, the actions are discrete; there are only two. Thus this corresponds to sampling actions from a Bernoulli distribution, where the probabilities of the left or right action come from the two softmax units.

```
ep_states.append(obs) #record the state
# pick a single random action for time step t, based on state obs
action_probs = sess.run(pi, feed_dict={x:[obs]})
#Get action sample using Bernoulli, probs already provided thanks to policy
5 if np.random.uniform(0,1) < action_probs[0][0]:
    action = 0
else:
    action = 1
10 ep_actions.append(action)
```



## Part 3

### (a) Printout:

```
At Episode 0 average number of time steps per episode = 14.0
Weights of the policy function:
[[ 0.55755877  0.60240793]
 [ 0.44874287 -0.53820491]
 [ 0.45429996 -0.63923603]
 [ 0.01130467  0.78500932]]
At Episode 2 average number of time steps per episode = 56.6666666667
Weights of the policy function:
[[ 0.53600371  0.623963  ]
 [ 0.43200225 -0.52146423]
 [ 0.45346969 -0.6384058  ]
 [ 0.00672296  0.78959101]]
At Episode 3 average number of time steps per episode = 51.0
Weights of the policy function:
[[ 0.53594983  0.62401688]
 [ 0.43245232 -0.5219143  ]
 [ 0.45365512 -0.63859123]
 [ 0.00674696  0.78956699]]
At Episode 100 average number of time steps per episode = 31.2376237624
Weights of the policy function:
[[ 0.53436333  0.62560338]
 [ 0.43284601 -0.52230793]
 [ 0.44868213 -0.63361824]
 [-0.01298162  0.80929571]]
At Episode 200 average number of time steps per episode = 34.0597014925
Weights of the policy function:
[[ 0.53912431  0.6208424  ]
 [ 0.46296576 -0.55242783]
 [ 0.44166663 -0.62660259]
 [-0.06721984  0.86353391]]
At Episode 300 average number of time steps per episode = 33.8239202658
Weights of the policy function:
[[ 0.53848225  0.62148446]
 [ 0.45933744 -0.54879963]
 [ 0.43758854 -0.62252426]
 [-0.08025246  0.87656689]]
At Episode 400 average number of time steps per episode = 34.6259351621
Weights of the policy function:
[[ 0.5386703  0.62129641]
 [ 0.45510232 -0.5445646  ]
 [ 0.43074602 -0.61568159]
 [-0.09911039  0.89542466]]
At Episode 500 average number of time steps per episode = 35.7884231537
Weights of the policy function:
[[ 0.53805286  0.62191385]
 [ 0.46096849 -0.5504303  ]
```

```
[ 0.41963291 -0.6045683 ]
[-0.14524838  0.94156253]]
At Episode 600 average number of time steps per episode = 36.3410981697
Weights of the policy function:
[[ 0.54066473  0.61930197]
 [ 0.43943968 -0.52890176]
 [ 0.41010731 -0.59504306]
 [-0.14948153  0.94579589]]
At Episode 700 average number of time steps per episode = 36.8045649073
Weights of the policy function:
[[ 0.53655952  0.62340719]
 [ 0.42556515 -0.5150274 ]
 [ 0.39745358 -0.58238965]
 [-0.1867042  0.98301834]]
At Episode 800 average number of time steps per episode = 37.7265917603
Weights of the policy function:
[[ 0.53741318  0.62255353]
 [ 0.43591556 -0.52537787]
 [ 0.38919657 -0.5741322 ]
 [-0.22438921  1.02070367]]
At Episode 900 average number of time steps per episode = 37.8213096559
Weights of the policy function:
[[ 0.53998774  0.61997896]
 [ 0.44810259 -0.53756505]
 [ 0.38740206 -0.57233781]
 [-0.23696616  1.03328073]]
At Episode 1000 average number of time steps per episode = 38.3526473526
Weights of the policy function:
[[ 0.54247153  0.61749518]
 [ 0.44651994 -0.53598231]
 [ 0.37451959 -0.55945563]
 [-0.2823396  1.07865429]]
At Episode 1100 average number of time steps per episode = 38.9554950045
Weights of the policy function:
[[ 0.54386896  0.61609775]
 [ 0.44097793 -0.53044027]
 [ 0.36343193 -0.5483678 ]
 [-0.31686878  1.11318338]]
At Episode 1200 average number of time steps per episode = 39.3896752706
Weights of the policy function:
[[ 0.5424608  0.61750591]
 [ 0.45023283 -0.53969496]
 [ 0.35200471 -0.53694052]
 [-0.37174591  1.16806066]]
At Episode 1300 average number of time steps per episode = 39.8239815527
Weights of the policy function:
[[ 0.54581171  0.61415499]
 [ 0.461882  -0.55134451]
 [ 0.34298766 -0.52792341]
```

```
[-0.42273465  1.21904969]]
At Episode 1400 average number of time steps per episode = 40.4354032834
Weights of the policy function:
[[ 0.54816878  0.61179787]
 [ 0.46878794 -0.55825049]
 [ 0.33294338 -0.51787919]
 [-0.46084109  1.25715744]]
At Episode 1500 average number of time steps per episode = 40.7868087941
Weights of the policy function:
[[ 0.55583572  0.60413092]
 [ 0.48687363 -0.57633632]
 [ 0.3293772  -0.5143128 ]
 [-0.48769423  1.28401065]]
At Episode 1600 average number of time steps per episode = 41.514678326
Weights of the policy function:
[[ 0.55211782  0.60784882]
 [ 0.4712818  -0.56074476]
 [ 0.31567404 -0.50060946]
 [-0.52434981  1.32066584]]
At Episode 1700 average number of time steps per episode = 41.8906525573
Weights of the policy function:
[[ 0.54059702  0.61936963]
 [ 0.45211741 -0.5415802 ]
 [ 0.30363086 -0.48856625]
 [-0.5473851  1.34370077]]
At Episode 1800 average number of time steps per episode = 42.3981121599
Weights of the policy function:
[[ 0.52419233  0.63577431]
 [ 0.43957868 -0.52904117]
 [ 0.29083389 -0.47576928]
 [-0.59951568  1.3958317 ]]
At Episode 1900 average number of time steps per episode = 42.87743293
Weights of the policy function:
[[ 0.51657081  0.64339584]
 [ 0.43919921 -0.52866197]
 [ 0.28493458 -0.46986997]
 [-0.62112528  1.41744184]]
At Episode 2000 average number of time steps per episode = 43.6636681659
Weights of the policy function:
[[ 0.5068773  0.65308934]
 [ 0.40527204 -0.49473462]
 [ 0.25522092 -0.44015631]
 [-0.70473766  1.50105453]]
At Episode 2100 average number of time steps per episode = 44.2793907663
Weights of the policy function:
[[ 0.51534802  0.64461863]
 [ 0.34820578 -0.43766838]
 [ 0.23294462 -0.41788   ]
 [-0.71354812  1.5098654 ]]
```

At Episode 2200 average number of time steps per episode = 45.0549750114

Weights of the policy function:

```
[[ 0.53165525  0.62831134]
 [ 0.31564504 -0.40510768]
 [ 0.21052435 -0.39545962]
 [-0.7434383   1.53975523]]
```

At Episode 2300 average number of time steps per episode = 45.5006518905

Weights of the policy function:

```
[[ 0.53189516  0.62807143]
 [ 0.30605838 -0.39552101]
 [ 0.20217578 -0.38711098]
 [-0.75433451  1.55065107]]
```

At Episode 2400 average number of time steps per episode = 45.9920866306

Weights of the policy function:

```
[[ 0.52435535  0.6356113 ]
 [ 0.2743797  -0.36384234]
 [ 0.18376607 -0.36870134]
 [-0.79971081  1.5960263 ]]
```

At Episode 2500 average number of time steps per episode = 46.419032387

Weights of the policy function:

```
[[ 0.55165166  0.60831499]
 [ 0.28457341 -0.37403607]
 [ 0.17400324 -0.35893857]
 [-0.82451653  1.62083161]]
```

At Episode 2600 average number of time steps per episode = 46.8904267589

Weights of the policy function:

```
[[ 0.55484182  0.60512483]
 [ 0.29693455 -0.38639712]
 [ 0.16573018 -0.35066542]
 [-0.88662434  1.68293941]]
```

At Episode 2700 average number of time steps per episode = 47.698630137

Weights of the policy function:

```
[[ 0.50808698  0.65187955]
 [ 0.20047094 -0.2899335 ]
 [ 0.13719249 -0.3221277 ]
 [-0.92898381  1.72529876]]
```

At Episode 2800 average number of time steps per episode = 48.3691538736

Weights of the policy function:

```
[[ 0.49883416  0.6611324 ]
 [ 0.18390864 -0.27337098]
 [ 0.11568674 -0.30062205]
 [-1.01573157  1.81204665]]
```

At Episode 2900 average number of time steps per episode = 48.9448466046

Weights of the policy function:

```
[[ 0.42597786  0.73398882]
 [ 0.11683179 -0.20629403]
 [ 0.09575897 -0.28069434]
 [-1.04848719  1.84480226]]
```

At Episode 3000 average number of time steps per episode = 49.7877374209

Weights of the policy function:

```
[[ 0.43547112  0.72449559]
 [ 0.12312732 -0.21258955]
 [ 0.0748591  -0.25979441]
 [-1.09224343  1.88855851]]
```

At Episode 3033 average number of time steps per episode = 50.0049439684

Weights of the policy function:

```
[[ 0.36846304  0.79150385]
 [ 0.04124827 -0.13071053]
 [ 0.06586583 -0.25080118]
 [-1.09515846  1.89147353]]
```

**(b)**

According to the setup codes below,

The input vector has a dimension of 4.  $(x, x', \theta, \theta')$

$x$  – position of cart along the horizontal axis, (with  $x=0$  being the middle point.  $x$  takes positive and negative values for the cart on different sides of the middle point)

$x'$  – velocity of cart along the horizontal axis

$\theta$  – angle between the pole and its balanced position (the pole is considered balanced if it is perpendicular to the horizontal axis)

$\theta'$  – angular velocity of the pole

Observing our output, we found that

```
def _step(self, action):
    assert self.action_space.contains(action), "%r (%s) invalid"%(action, type(action))
    state = self.state
    x, x_dot, theta, theta_dot = state
    force = self.force_mag if action==1 else -self.force_mag
    costheta = math.cos(theta)
    sintheta = math.sin(theta)
    temp = (force + self.polemass_length * theta_dot * theta_dot * sintheta) /
           self.total_mass
    thetaacc = (self.gravity * sintheta - costheta*temp) / (self.length *
        (4.0/3.0 - self.masspole * costheta * costheta / self.total_mass))
    xacc = temp - self.polemass_length * thetaacc * costheta / self.total_mass
    x = x + self.tau * x_dot
    x_dot = x_dot + self.tau * xacc
    theta = theta + self.tau * theta_dot
    theta_dot = theta_dot + self.tau * thetaacc
    self.state = (x,x_dot,theta,theta_dot)
    done = x < -self.x_threshold \
           or x > self.x_threshold \
           or theta < -self.theta_threshold_radians \
           or theta > self.theta_threshold_radians
    done = bool(done)

    return np.array(self.state), reward, done, {}
```