

Mastering Atari with an Actor-Critic Model

Sam Greydanus
greydasa@oregonstate.edu

January 19, 2018

Abstract

Reinforcement Learning (RL) with deep neural networks is an exciting area of research, but training a deep RL agent is difficult in practice. In this talk, I will introduce the theories and intuitions of the (Asynchronous) Advantage Actor-Critic (A3C) model and walk through a 180-line implementation in PyTorch. By the end of my talk, you should have the tools you need to train your own Atari agents.

1 A3C in theory

The idea behind reinforcement learning is that, given an *agent* which can interact with its *environment*, one can learn a policy that maximizes the expected reward.

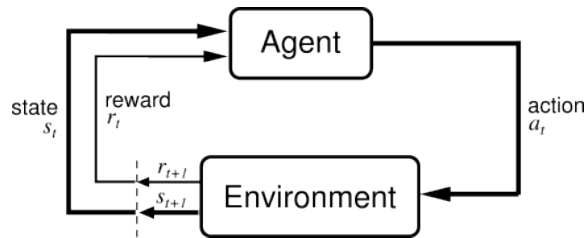


Figure 1: Overview of the RL framework. An agent in state s_t receives reward r_t from the environment. Then, it chooses some internal policy, π to choose action a_t , interacts with the environment, and receives the next state/reward pair.

The idea is simple, but making it work in practice is tricky for a variety of reasons. In this talk, I'll use the Atari video games as example environments because they show off 1) why getting RL to work in practice is tricky 2) why RL is a powerful and interesting tool.

1.1 Policy Gradient Theorem

Suppose we have a reward function $f(x)$ and an agent that acts according to probability distribution $p(x)$. We would like to modify the distribution $p(x)$ to maximize the expectation value $E_x[f(x)]$. Let's assume $p(x)$ is differentiable and we can optimize it using gradient descent with respect to the expected reward.

Now all we need is to take the gradient of the expectation value:

$$\nabla_{\theta} E_x[f(x)] = \nabla_{\theta} \sum_x p(x) f(x) \quad \text{definition of expectation value} \quad (1)$$

$$= \sum_x \nabla_{\theta} p(x) f(x) \quad \text{apply gradient to each term} \quad (2)$$

$$= \sum_x p(x) \frac{\nabla_{\theta} p(x)}{p(x)} f(x) \quad \text{multiply by identity} \quad (3)$$

$$= \sum_x p(x) \nabla_{\theta} \log p(x) f(x) \quad \text{because } \nabla_{\theta} \log z = \frac{1}{z} \nabla_{\theta} z \quad (4)$$

$$= E_x[\nabla_{\theta} \log p(x) f(x)] \quad \text{this is the Policy Gradient Theorem} \quad (5)$$

This gives us the gradient of the expected reward. From now on, let's start using π_{θ} as our $p(x)$ and call it the agent's *policy*. We'll also refer to the expected reward. This gives us the gradient of the expected reward. From now on, let's start using π_{θ} as our $p(x)$ and call it the agent's *policy*. We'll also refer to the expected reward $E_x[f(x)]$ as J .

1.2 REINFORCE

Many times, our agent will receive a reward not only because of the action a_t which it just took, but also because of the actions which led it into the high-reward state. These actions consist of the set a_{t-1}, a_{t-2}, \dots . We would like to give the agent credit for these actions, especially the ones which helped it obtain the final reward.

There is more than one solution to this problem of *credit assignment*, but here's one idea. As you get closer to the high-reward state, the actions you take are exponentially more likely to have contributed to that reward. If this is the case, then you should spread the reward you received backwards in time. We do this via *gamma-discounted rewards*:

$$R = \sum_{l=0}^{\infty} \gamma^l r_{t+l} \quad (6)$$

Here, γ is a hyperparameter that we, as experimenters, get to set. It generally corresponds to how sparse the rewards are. We will use $\gamma = 0.99$ for all experiments. We'll use the value function $V^{\pi}(s) = E_s[R]$ to represent the value of the state. If we want to describe the value of taking a particular action once we're in this state, we use the Q function, $Q^{\pi}(s, a) = R(s, a) + \gamma V^*(\delta(s, a))$.

Now our policy gradient looks like this:

$$\nabla J(\theta) = E_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi}(s, a)] \quad (7)$$

Notice that this is just the gradient of a differentiable function, $\log \pi_{\theta}(s, a)$, times a scalar which we know how to compute, $Q^{\pi}(s, a)$. If we want to, we can go ahead and train an agent using gradient descent. This, in fact, is an algorithm called REINFORCE and it works well in some simple cases. Sadly, it doesn't work great for Atari.

Note. Many times, people will clip raw rewards to $[-1, 1]$ before gamma discounting. I tried training on Atari without this step and all my models failed.

1.3 Actor-critic

The big remaining issue is that $Q^\pi(s, a)$ is not a well-behaved function (especially for Atari and other complex environments). It is sparse, high-variance, and totally unnormalized.

The actor-critic trick is meant to reduce the variance of the policy gradient. Instead of using $Q^\pi(s, a)$ directly, we learn an estimator,

$$Q_w(s, a) \approx Q^\pi(s, a) \quad (8)$$

This estimate has much lower variance, leading to better results in practice. Now we have:

$$\nabla J(\theta) \approx E_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)] \quad (9)$$

Question. Why not simply reduce variance by keeping a running estimate of batch statistics and then using them to reduce variance? Wouldn't this accomplish the same "smoothing" objective?

Answer. I'm not sure. I tried doing this in practice and it did not work as well. I've seen other people do this *within* batches, which seemed to help with REINFORCE.

1.4 Advantage actor-critic

Another way to reduce the variance is to subtract a baseline from the policy gradient. To see how this works, imagine if our expected reward looked like this:

$$E_{\pi_\theta} [\log \pi_\theta(s, a) (f(s, a) - g(s))] \quad (10)$$

I want to convince you that if g has zero mean, then it won't change the expectation value:

$$E_{\pi_\theta} [\log \pi_\theta(s, a) g(s)] = \sum_{s \in S} \sum_{a \in A} \pi_\theta(s, a) g(s) \quad \text{from eqn. 2} \quad (11)$$

$$= \sum_{s \in S} g(s) \sum_{a \in A} \pi_\theta(s, a) \quad (12)$$

$$= \sum_{s \in S} g(s) \quad \text{because } \sum_{a \in A} \pi_\theta(s, a) = 1 \quad (13)$$

$$= 0 \quad \text{if } g \text{ has zero mean} \quad (14)$$

Now, because $E[A - B] = E[A] - E[B]$, we can conclude that

$$E_{\pi_\theta} [\log \pi_\theta(s, a) (f(s, a) - g(s))] = E_{\pi_\theta} [\log \pi_\theta(s, a) f(s, a)] \quad (15)$$

People call this "subtracting a baseline" and it is another way to reduce variance. A good baseline ends up being the value function, $V^\pi(s)$. Recall that the advantage function is defined as

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) \quad (16)$$

$$\approx Q_w(s, a) - V_v(s) \quad (17)$$

In equation 17 I've replaced the value function with a learned estimate, just as we did for the Q -values in section 1.3. So, our new policy gradient becomes

$$\nabla J(\theta) \approx E_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^\pi(s, a)] \quad (18)$$

1.5 TD actor-critic

One remaining issue is that, if we want to estimate the advantage using a neural network, we need two separate estimators, one for $Q_w(s, a)$ and one for $V_v(s)$. There's a better way: using the TD error, we can estimate the advantage directly from $V_v(s)$. Why? The TD error is an unbiased estimate of the advantage function. Recall that the TD error is given by

$$\delta^{\pi_\theta} = r + \gamma V^\pi(s_{t+1}) - V^\pi(s) \quad (19)$$

The expectation value of the TD error looks like:

$$E[\delta^{\pi_\theta} | s, a] = E[r + \gamma V^\pi(s_{t+1}) | s, a] - V^\pi(s) \quad (20)$$

$$= Q^\pi(s, a) - V^\pi(s) \quad (21)$$

$$= A^\pi(s, a) \quad (22)$$

Thus, if we use TD error, we only need to train one value estimator to approximate the advantage. The policy gradient for our final advantage actor-critic model looks like:

$$\nabla J(\theta) \approx E_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) \delta^{\pi_\theta}] \quad (23)$$

This approximation produces gradients with much lower variance, giving a deep neural network enough signal to learn to play Atari games very well.

2 A3C in practice

Making deep RL work in practice an art form unto itself. It takes a lot of clever engineering! As a result, most open-source implementations are large, clunky, and confusing. Otherwise, they are small, overly simplistic, and can't handle Atari. Frustrated by this, I ended up writing my own version called `baby-a3c`¹ which makes solving the Atari environments as simple as possible...but not simpler!

I've appended snapshots of the code below.

¹github.com/greydanus/baby-a3c

```

1 # Baby Advantage Actor-Critic | Sam Greydanus | October 2017 | MIT License
2
3 from __future__ import print_function
4 import torch, os, gym, time, glob, argparse
5 import numpy as np
6 from scipy.signal import lfilter
7 from scipy.misc import imresize # preserves single-pixel info _unlike_ img = img[:, :, 2]
8
9 import torch.nn as nn
10 from torch.autograd import Variable
11 import torch.nn.functional as F
12 import torch.multiprocessing as mp
13 os.environ['OMP_NUM_THREADS'] = '1'
14
15 parser = argparse.ArgumentParser(description=None)
16 parser.add_argument('--env', default='Breakout-v0', type=str, help='gym environment')
17 parser.add_argument('--processes', default=20, type=int, help='number of processes to train with')
18 parser.add_argument('--render', default=False, type=bool, help='renders the atari environment')
19 parser.add_argument('--test', default=False, type=bool, help='test mode sets lr=0, chooses most likely actions')
20 parser.add_argument('--lstm_steps', default=20, type=int, help='steps to train LSTM over')
21 parser.add_argument('--lr', default=1e-4, type=float, help='learning rate')
22 parser.add_argument('--seed', default=1, type=int, help='seed random # generators (for reproducibility)')
23 parser.add_argument('--gamma', default=0.99, type=float, help='discount for gamma-discounted rewards')
24 parser.add_argument('--tau', default=1.0, type=float, help='discount for generalized advantage estimation')
25 parser.add_argument('--horizon', default=0.99, type=float, help='horizon for running averages')
26 args = parser.parse_args()
27
28 args.save_dir = '{}/'.format(args.env.lower()) # keep the directory structure simple
29 if args.render: args.processes = 1 ; args.test = True # render mode -> test mode w one process
30 if args.test: args.lr = 0 # don't train in render mode
31 args.num_actions = gym.make(args.env).action_space.n # get the action space of this game
32 os.makedirs(args.save_dir) if not os.path.exists(args.save_dir) else None # make dir to save models etc.
33
34 discount = lambda x, gamma: lfilter([1],[1,-gamma],x[::-1])[::-1] # discounted rewards one liner
35 prepro = lambda img: imresize(img[35:195].mean(2), (80,80)).astype(np.float32).reshape(1,80,80)/255.
36
37 def printlog(args, s, end='\n', mode='a'):
38     print(s, end=end) ; f=open(args.save_dir+'log.txt',mode) ; f.write(s+'\n') ; f.close()
39
40 class NNPolicy(torch.nn.Module): # an actor-critic neural network
41     def __init__(self, channels, num_actions):
42         super(NNPolicy, self).__init__()
43         self.conv1 = nn.Conv2d(channels, 32, 3, stride=2, padding=1)
44         self.conv2 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
45         self.conv3 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
46         self.conv4 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
47         self.lstm = nn.LSTMCell(32 * 5 * 5, 256)
48         self.critic_linear, self.actor_linear = nn.Linear(256, 1), nn.Linear(256, num_actions)
49
50     def forward(self, inputs):
51         inputs, (hx, cx) = inputs
52         x = F.elu(self.conv1(inputs))
53         x = F.elu(self.conv2(x))
54         x = F.elu(self.conv3(x))
55         x = F.elu(self.conv4(x))
56         hx, cx = self.lstm(x.view(-1, 32 * 5 * 5), (hx, cx))
57         return self.critic_linear(hx), self.actor_linear(hx), (hx, cx)
58
59     def try_load(self, save_dir):
60         paths = glob.glob(save_dir + '*.tar') ; step = 0
61         if len(paths) > 0:
62             ckpts = [int(s.split('.')[2]) for s in paths]
63             ix = np.argmax(ckpts) ; step = ckpts[ix]
64             self.load_state_dict(torch.load(paths[ix]))
65         print("\tno saved models") if step is 0 else print("\tloaded model: {}".format(paths[ix]))
66         return step

```

```

68 class SharedAdam(torch.optim.Adam): # extend a pytorch optimizer so it shares grads across processes
69     def __init__(self, params, lr=1e-3, betas=(0.9, 0.999), eps=1e-8, weight_decay=0):
70         super(SharedAdam, self).__init__(params, lr, betas, eps, weight_decay)
71         for group in self.param_groups:
72             for p in group['params']:
73                 state = self.state[p]
74                 state['shared_steps'], state['step'] = torch.zeros(1).share_memory_(), 0
75                 state['exp_avg'] = p.data.new().resize_as_(p.data).zero_().share_memory_()
76                 state['exp_avg_sq'] = p.data.new().resize_as_(p.data).zero_().share_memory_()
77
78         def step(self, closure=None):
79             for group in self.param_groups:
80                 for p in group['params']:
81                     if p.grad is None: continue
82                     self.state[p]['shared_steps'] += 1
83                     self.state[p]['step'] = self.state[p]['shared_steps'][0] - 1 # there's a "step += 1" later
84             super.step(closure)
85
86 torch.manual_seed(args.seed)
87 shared_model = NNPolicy(channels=1, num_actions=args.num_actions).share_memory()
88 shared_optimizer = SharedAdam(shared_model.parameters(), lr=args.lr)
89
90 info = {k : torch.DoubleTensor([0]).share_memory_() for k in ['run_epr', 'run_loss', 'episodes', 'frames']}
91 info['frames'] += shared_model.try_load(args.save_dir)*1e6
92 if int(info['frames'][0]) == 0: printlog(args, '', end='', mode='w') # clear log file
93
94 def cost_func(values, logps, actions, rewards):
95     np_values = values.view(-1).data.numpy()
96
97     # generalized advantage estimation (a policy gradient method)
98     delta_t = np.asarray(rewards) + args.gamma * np_values[1:] - np_values[:-1]
99     gae = discount(delta_t, args.gamma * args.tau)
100     logpys = logps.gather(1, Variable(actions).view(-1,1))
101     policy_loss = -(logpys.view(-1) * Variable(torch.Tensor(gae))).sum()
102
103     # l2 loss over value estimator
104     rewards[-1] += args.gamma * np_values[-1]
105     discounted_r = discount(np.asarray(rewards), args.gamma)
106     discounted_r = Variable(torch.Tensor(discounted_r))
107     value_loss = .5 * (discounted_r - values[:-1,0]).pow(2).sum()
108
109     entropy_loss = -(-logps * torch.exp(logps)).sum() # encourage lower entropy
110     return policy_loss + 0.5 * value_loss + 0.01 * entropy_loss

```

```

112 def train(rank, args, info):
113     env = gym.make(args.env) # make a local (unshared) environment
114     env.seed(args.seed + rank) ; torch.manual_seed(args.seed + rank) # seed everything
115     model = NNPolicy(channels=1, num_actions=args.num_actions) # init a local (unshared) model
116     state = torch.Tensor(prepro(env.reset())) # get first state
117
118     start_time = last_disp_time = time.time()
119     episode_length, epr, eploss, done = 0, 0, 0, True # bookkeeping
120
121     while info['frames'][0] <= 8e7 or args.test: # openai baselines uses 40M frames...we'll use 80M
122         model.load_state_dict(shared_model.state_dict()) # sync with shared model
123
124         cx = Variable(torch.zeros(1, 256)) if done else Variable(cx.data) # lstm memory vector
125         hx = Variable(torch.zeros(1, 256)) if done else Variable(hx.data) # lstm activation vector
126         values, logps, actions, rewards = [], [], [], [] # save values for computing gradients
127
128         for step in range(args.lstm_steps):
129             episode_length += 1
130             value, logit, (hx, cx) = model((Variable(state.view(1,1,80,80)), (hx, cx)))
131             logp = F.log_softmax(logit)
132
133             action = logp.max(1)[1].data if args.test else torch.exp(logp).multinomial().data[0]
134             state, reward, done, _ = env.step(action.numpy()[0])
135             if args.render: env.render()
136
137             state = torch.Tensor(prepro(state)) ; epr += reward
138             reward = np.clip(reward, -1, 1) # reward
139             done = done or episode_length >= 1e4 # keep agent from playing one episode too long
140
141             info['frames'] += 1 ; num_frames = int(info['frames'][0])
142             if num_frames % 2e6 == 0: # save every 2M frames
143                 printlog(args, '\n\t{:.0f}M frames: saved model\n'.format(num_frames/1e6))
144                 torch.save(shared_model.state_dict(), args.save_dir+'model_{:.0f}.tar'.format(num_frames/1e6))
145
146             if done: # update shared data. maybe print info.
147                 info['episodes'] += 1
148                 interp = 1 if info['episodes'][0] == 1 else 1 - args.horizon
149                 info['run_epr'].mul_(1-interp).add_(interp * epr)
150                 info['run_loss'].mul_(1-interp).add_(interp * eploss)
151
152                 if rank == 0 and time.time() - last_disp_time > 60: # print info ~ every minute
153                     elapsed = time.strftime("%Hh %Mm %Ss", time.gmtime(time.time() - start_time))
154                     printlog(args, 'time {}, episodes {:.0f}, frames {:.1f}M, run epr {:.2f}, run loss {:.2f}'
155                             .format(elapsed, info['episodes'][0], num_frames/1e6, info['run_epr'][0], info['run_loss'][0]))
156                     last_disp_time = time.time()
157
158                 episode_length, epr, eploss = 0, 0, 0
159                 state = torch.Tensor(prepro(env.reset()))
160
161             values.append(value) ; logps.append(logp) ; actions.append(action) ; rewards.append(reward)
162
163             next_value = Variable(torch.zeros(1,1)) if done else model((Variable(state.unsqueeze(0)), (hx, cx)))[0]
164             values.append(Variable(next_value.data))
165
166             loss = cost_func(torch.cat(values), torch.cat(logps), torch.cat(actions), np.asarray(rewards))
167             eploss += loss.data[0]
168             shared_optimizer.zero_grad() ; loss.backward()
169             torch.nn.utils.clip_grad_norm(model.parameters(), 40)
170
171             for param, shared_param in zip(model.parameters(), shared_model.parameters()):
172                 if shared_param.grad is None: shared_param._grad = param.grad # sync gradients with shared model
173             shared_optimizer.step()
174
175 processes = []
176 for rank in range(args.processes):
177     p = mp.Process(target=train, args=(rank, args, info))
178     p.start() ; processes.append(p)
179 for p in processes:
180     p.join()

```