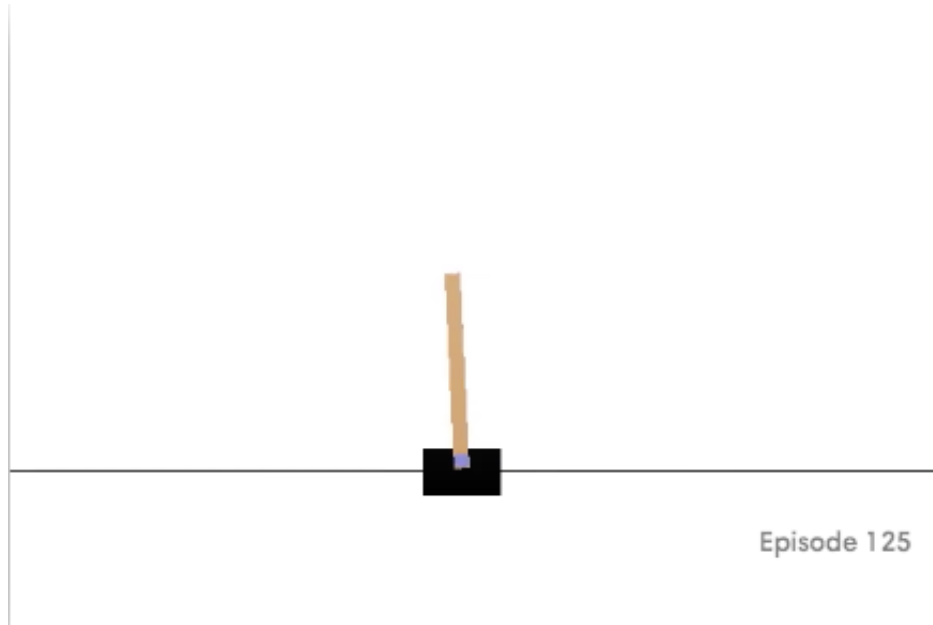


# Learn to Play CartPole with PyTorch using PGPE

---

Code: [agent.py](#), [cartpole.py](#)

The CartPole environment:



"Keep the pole from falling for a short period of time", that is the task our algorithm will learn to perform in this post.

See the [final result](#) on my OpenAI Gym page

See the [sample code](#) on my Github page

## INTRODUCTION

In this post, we will apply another policy-based algorithm to solve the CartPole environment. Unlike REINFORCE, searching the optimal policy in the action space, policy gradient with parameter-based explorations (PGPE) searches the optimal policy in the parameter space. The idea is to equip each model parameter with a *prior distribution* to sample from. At the beginning of each roll-out, we sample a set of model parameters, and follow a *deterministic* controller (policy) for the episode thereafter. Just like REINFORCE, the original PGPE adopts an episodic setting. After an episode finishes, we compute the gradient of the objective function w.r.t. *hyper-parameters of the prior distribution* and do the parameter update.

## ALGORITHM

Note:

- Typing LaTeX inline is not supported by my Markdown editor, so some of the equations may be a little difficult to read, refer to the [original paper](#) for details.

- Some notations have been modified to better fit in the context, hence you may find some inconsistency in notations with the original paper.

1. Each action follows the *deterministic policy* (the  $\delta$  term) with model parameters  $\delta$  drawn from the *prior distribution*, parametrized by  $\rho$

$$p(a|s, \rho) = \int_{\Theta} p(\theta|\rho) \delta_{F_{\theta}(s)=a} d\theta$$

which means the probability of taking action  $a$  under state  $s$  is exactly the probability of drawing parameter  $\theta$  given hyper-parameter  $\rho$ , since the delta function (our deterministic controller) contributes no stochasticity in the equation.

2. The objective  $J(\rho)$  is defined to be

$$J(\rho) = \int_{\Theta} \int_H p(h, \theta|\rho) r(h) dh d\theta$$

which is the *expected return* taken over all possible parameters we can draw from the prior and the history generated by following policy  $p(h, \theta)$ .

3. Apply the classic log-likelihood trick, we obtain the gradient of the objective w.r.t. hyper-parameter  $\rho$

$$\nabla_{\rho} J(\rho) = \int_{\Theta} \int_H p(h, \theta|\rho) \nabla_{\rho} \log p(h, \theta|\rho) r(h) dh d\theta$$

4. Observe that the history is conditionally independent of  $\rho$  given  $\theta$ , that is,  $p(h, \theta|\rho) = p(h|\rho)p(\theta|\rho)$  and the partial derivative reduces,  $\nabla_{\rho} \log p(h, \theta|\rho) = \nabla_{\rho} \log p(\theta|\rho)$ . Combining all of the above, the gradient estimator reveals:

$$\nabla_{\rho} J(\rho) \approx \frac{1}{N} \sum_{n=1}^N \nabla_{\rho} \log p(\theta^n|\rho) r(h^n)$$

**Remark** the above gradient estimator is completely determined by the agent's model parameters *without* the knowledge of environment dynamics. Information concerning environment dynamics encapsulates in the term  $p(h|\theta)$ , as we get rid of it when taking the log value of the term and differentiate w.r.t.  $\rho$ . Hence PGPE is a *model-free* algorithm, just as REINFORCE.

5. Here we assume each model parameter  $\theta_i$  follows a normal distribution independent with one another. Hence the hyper-parameters  $\rho$  consists of a sequence of means and standard deviations  $\rho = ((\mu_i, \sigma_i)_i)$ . Direct computation gives

$$\nabla_{\mu} \log p(\theta|\rho) = \frac{\theta - \mu}{\sigma^2}, \quad \nabla_{\sigma} \log p(\theta|\rho) = \frac{(\theta - \mu)^2 - \sigma^2}{\sigma^3}$$

Finally, we choose a stepsize proportional to the variance,  $\alpha_{\rho} = \alpha \sigma^2$ . The stepsize may be different for different parameters as we will see below.

$$\Delta \mu = \alpha_{\mu} (r - b) (\theta - \mu), \quad \Delta \sigma = \alpha_{\sigma} (r - b) \frac{(\theta - \mu)^2 - \sigma^2}{\sigma}$$

where  $b$  is the reward baseline for variance reduction purpose.

**Algorithm 1** The PGPE Algorithm without reward normalization: Left side shows the basic version, right side shows the version with symmetric sampling.  $\mathbf{T}$  and  $\mathbf{S}$  are  $P \times N$  matrices with  $P$  the number of parameters and  $N$  the number of histories. The baseline is updated accordingly after each step.  $\alpha$  is the learning rate or step size.

|  |  |
|--|--|
| Initialize $\mu$ to $\mu_{\text{init}}$  | Initialize $\mu$ to $\mu_{\text{init}}$  |
| Initialize $\sigma$ to $\sigma_{\text{init}}$  | Initialize $\sigma$ to $\sigma_{\text{init}}$  |
| <b>while</b> TRUE <b>do</b>  | <b>while</b> TRUE <b>do</b>  |
| <b>for</b> $n = 1$ to $N$ <b>do</b>  | <b>for</b> $n = 1$ to $N$ <b>do</b>  |
| draw $\theta^n \sim \mathcal{N}(\mu, I\sigma^2)$                                     | draw perturbation $\epsilon^n \sim \mathcal{N}(0, I\sigma^2)$                                      |
|  | $\theta^{+,n} = \mu + \epsilon^n$  |
|  | $\theta^{-,n} = \mu - \epsilon^n$  |
| evaluate $r^n = r(h(\theta^n))$  | evaluate $r^{+,n} = r(h(\theta^{+,n}))$  |
|  | evaluate $r^{-,n} = r(h(\theta^{-,n}))$  |
| <b>end for</b>   | <b>end for</b>   |
| $\mathbf{T} = [t_{ij}]_{ij}$ with $t_{ij} := (\theta_i^j - \mu_i)$                   | $\mathbf{T} = [t_{ij}]_{ij}$ with $t_{ij} := \epsilon_i^j$   |
| $\mathbf{S} = [s_{ij}]_{ij}$ with $s_{ij} := \frac{t_{ij}^2 - \sigma_i^2}{\sigma_i}$ | $\mathbf{S} = [s_{ij}]_{ij}$ with $s_{ij} := \frac{(\epsilon_i^j)^2 - \sigma_i^2}{\sigma_i}$       |
| $\mathbf{r} = [(r^1 - b), \dots, (r^N - b)]^T$                                       | $\mathbf{r}_T = [(r^{+,1} - r^{-,1}), \dots, (r^{+,N} - r^{-,N})]^T$                               |
|  | $\mathbf{r}_S = [\frac{(r^{+,1} + r^{-,1})}{2} - b, \dots, (\frac{(r^{+,N} + r^{-,N})}{2} - b)]^T$ |
| update $\mu = \mu + \alpha \mathbf{T} \mathbf{r}$                                    | update $\mu = \mu + \alpha \mathbf{T} \mathbf{r}_T$  |
| update $\sigma = \sigma + \alpha \mathbf{S} \mathbf{r}$                              | update $\sigma = \sigma + \alpha \mathbf{S} \mathbf{r}_S$  |
| update baseline $b$ accordingly  | update baseline $b$ accordingly  |
| <b>end while</b>   | <b>end while</b>   |

The algorithm taken from the [original paper](#) is shown above. In this post we will implement the one on the left-hand side (the vanilla version). But as we will see later, this simple algorithm turned out to be very effective against benchmark control problems compared with to standard PG methods such as REINFORCE.

## CODE

*cartpole.py*

```

1  import os.path
2  import sys
3  import gym
4  from gym import wrappers
5  from agent import PGPE
6
7  env = gym.make('CartPole-v0')
8
9  # Gym's built-in monitor functionality
10 # We can later upload our result to OpenAI Gym
11 mod_path = os.path.dirname(os.path.abspath(sys.argv[0]))
12 save_path = os.path.join(mod_path, 'cartpole_experiment_1')
13
14 env = wrappers.Monitor(env, save_path, force=True)
15
16 # Create an agent instance
17 RL = PGPE(
18     n_features=env.observation_space.shape[0],
19     n_actions=env.action_space.n
20 )
21
22
23 for ep in range(200):
24     observation = env.reset()
25
26     while True:
27         # Typical RL env-agent paradigm
28         action = RL.choose_action(observation)
29         observation_, reward, done, info = env.step(action)
30
31         RL.store_reward(reward)
32         r = RL.get_return()
33
34         if done:
35             print("Episode:", ep, "   Reward:", int(r))
36             vt = RL.learn_and_sample()
37             break
38
39         observation = observation_
40
41 # Close env.
42 env.close()

```

*agent.py*

```

1  import numpy as np
2

```

```

3 import torch
4 import torch.nn as nn
5 from torch.autograd import Variable
6
7
8 class PGPE:
9     def __init__(self, n_actions, n_features):
10         self.n_actions = n_actions
11         self.n_features = n_features
12
13         # Return tracker for episode
14         self.ret = 0.
15
16         self.model = nn.Sequential(
17             nn.Linear(self.n_features, self.n_actions, bias=False),
18             nn.Softmax() # deterministic policy, pick action with
greater value
19         )
20
21         # Learning rate for hyper-params mu and sigma
22         self.Mu_lr = 0.2
23         self.Sigma_lr = 0.1
24
25         # Prepare hyper-params, store mean/var separately in lists
26         self.Param = list(self.model.parameters())
27         self.Mu = []
28         self.Sigma = []
29         for p in self.Param:
30             # initialize hyper-params
31             self.Mu.append(torch.normal(torch.zeros(p.size()),
torch.ones(p.size()))))
32             self.Sigma.append(2 * torch.ones(p.size()))
33
34             # Sample initial model params
35             p.data = torch.normal(self.Mu[-1], self.Sigma[-1])
36
37     def choose_action(self, obs):
38         # Scale input to (-1, 1):
39         # 1. if range is finite -> divide by range
40         # 2. if range is infinite -> take tanh
41         obs[0] /= 2.4
42         obs[1] = np.tanh(obs[1])
43         obs[2] /= 41.8
44         obs[3] = np.tanh(obs[3])
45
46         # cast np array to torch variable
47         s =
Variable(torch.from_numpy(obs.astype(np.float32))).unsqueeze(0)
48

```

```

49         a = self.model.forward(s).data.numpy()
50         action = a[0].argmax() # pick action with greater value
51         return action
52
53     def get_return(self):
54         return self.ret
55
56     def store_reward(self, r):
57         self.ret += r # Compute (un-discounted) return
58
59     def learn_and_sample(self):
60         # Scale return to [0, 1]
61         _r = self.ret / 200
62
63         # reset return tracker
64         self.ret = 0.
65
66         for i in range(len(self.Param)):
67             # Learning
68             # These are the T and S matrices in the original paper
69             _T = self.Param[i].data - self.Mu[i]
70             _S = (_T ** 2 - self.Sigma[i] ** 2) / self.Sigma[i]
71
72             # Update means
73             _delta_Mu = self.Mu_lr * _r * _T
74             self.Mu[i] += _delta_Mu
75
76             # Update standard deviations
77             _delta_Sigma = self.Sigma_lr * _r * _S
78             self.Sigma[i] += _delta_Sigma
79
80             # Freeze params if hit target reward, else re-sample
81             if _r < 1.:
82                 self.Param[i].data = torch.normal(self.Mu[i],
self.Sigma[i])

```

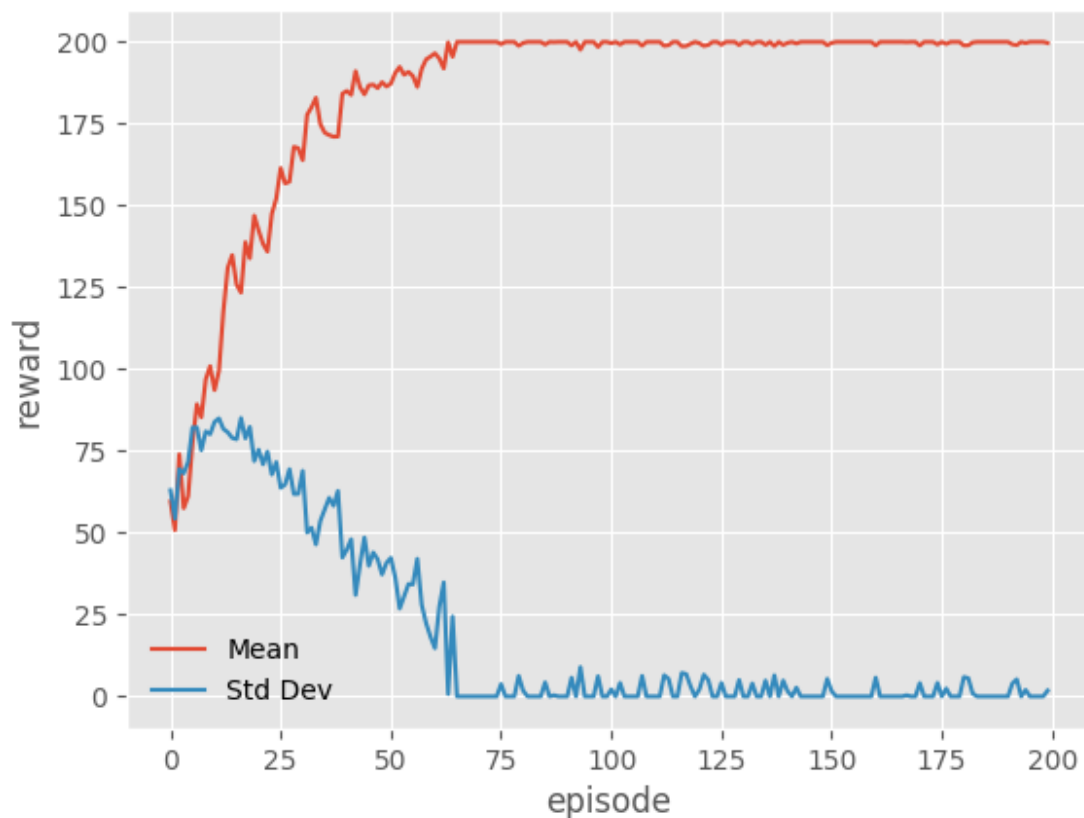
## RESULT

The result of this algorithm can be viewed on [my OpenAI Gym page](#). I also provide a gist link there so you can easily reproduce the result. As you can see, the algorithm converges quickly, taking only 5 episodes to solve the CartPole environment. The model consists of only 8 parameters and without performing any sophisticated initialization trick (in fact, random initialization from  $\text{normal}(0, 2)$  is used to draw initial model parameters). Other results provided on the [OpenAI Gym cartpole environment page](#) utilizing the REINFORCE algo often take up dozens more episodes to converge to a stable policy.

By the way, in case anyone is interested in what the input states are in the cartpole env. Here is the [official OpenAI gym github page](#) providing the background information. Note the env setting is a bit different from the one on the wiki page. But it doesn't really matter since we scale the input state when it comes in.

## ROBUSTNESS

To further evaluate our algorithm, we roll out 30 trials, each with 200 episodes. Then compute the mean and standard deviation of reward *at each episode* over these 30 trials to get a sense of how much reward in average we can get for each episode. The plot demonstrates the result:



To interpret, say, at the 50th episode, the algorithm gets an average of 187.5 and a standard deviation of reward of 37.5. The plot demonstrates the desired result: mean reward grows steadily toward 200 and the standard deviation decreases to 0 as the algorithm learns to prevent the pole from falling down for 200 steps. Also, we can see the algorithm *solves* the env. at around 63th step, where it achieves an average reward of 195.0 over 100 consecutive episodes thereafter.