

Part 3: Intro to Policy Optimization

Table of Contents

- Part 3: Intro to Policy Optimization
 - Deriving the Simplest Policy Gradient
 - Implementing the Simplest Policy Gradient
 - Expected Grad-Log-Prob Lemma
 - Don't Let the Past Distract You
 - Implementing Reward-to-Go Policy Gradient
 - Baselines in Policy Gradients
 - Other Forms of the Policy Gradient
 - Recap

In this section, we'll discuss the mathematical foundations of policy optimization algorithms, and connect the material to sample code. We will cover three key results in the theory of **policy gradients**:

- **the simplest equation** describing the gradient of policy performance with respect to policy parameters,
- a rule which allows us to **drop useless terms** from that expression,
- and a rule which allows us to **add useful terms** to that expression.

In the end, we'll tie those results together and describe the advantage-based expression for the policy gradient—the version we use in our [Vanilla Policy Gradient](#) implementation.

Deriving the Simplest Policy Gradient

Here, we consider the case of a stochastic, parameterized policy, π_θ . We aim to maximize the expected return $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$. For the purposes of this derivation, we'll take $R(\tau)$ to give the **finite-horizon undiscounted return**, but the derivation for the infinite-horizon discounted return setting is almost identical.

We would like to optimize the policy by gradient ascent, eg

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_\theta)|_{\theta_k}.$$

The gradient of policy performance, $\nabla_\theta J(\pi_\theta)$, is called the **policy gradient**, and algorithms that optimize the policy this way are called **policy gradient algorithms**. (Examples include Vanilla Policy Gradient and TRPO. PPO is often referred to as a policy gradient algorithm,

though this is slightly inaccurate.)

To actually use this algorithm, we need an expression for the policy gradient which we can numerically compute. This involves two steps: 1) deriving the analytical gradient of policy performance, which turns out to have the form of an expected value, and then 2) forming a sample estimate of that expected value, which can be computed with data from a finite number of agent-environment interaction steps.

In this subsection, we'll find the simplest form of that expression. In later subsections, we'll show how to improve on the simplest form to get the version we actually use in standard policy gradient implementations.

We'll begin by laying out a few facts which are useful for deriving the analytical gradient.

1. Probability of a Trajectory. The probability of a trajectory $\tau = (s_0, a_0, \dots, s_{T+1})$ given that actions come from π_θ is

$$P(\tau|\theta) = \rho_0(s_0) \prod_{t=0}^T P(s_{t+1}|s_t, a_t) \pi_\theta(a_t|s_t).$$

2. The Log-Derivative Trick. The log-derivative trick is based on a simple rule from calculus: the derivative of $\log x$ with respect to x is $1/x$. When rearranged and combined with chain rule, we get:

$$\nabla_\theta P(\tau|\theta) = P(\tau|\theta) \nabla_\theta \log P(\tau|\theta).$$

3. Log-Probability of a Trajectory. The log-prob of a trajectory is just

$$\log P(\tau|\theta) = \log \rho_0(s_0) + \sum_{t=0}^T \left(\log P(s_{t+1}|s_t, a_t) + \log \pi_\theta(a_t|s_t) \right).$$

4. Gradients of Environment Functions. The environment has no dependence on θ , so gradients of $\rho_0(s_0)$, $P(s_{t+1}|s_t, a_t)$, and $R(\tau)$ are zero.

5. Grad-Log-Prob of a Trajectory. The gradient of the log-prob of a trajectory is thus

$$\begin{aligned} \nabla_\theta \log P(\tau|\theta) &= \cancel{\nabla_\theta \log \rho_0(s_0)} + \sum_{t=0}^T \left(\cancel{\nabla_\theta \log P(s_{t+1}|s_t, a_t)} + \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \\ &= \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t). \end{aligned}$$

Putting it all together, we derive the following:

⚠ Derivation for Basic Policy Gradient

$$\begin{aligned}
 \nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] \\
 &= \nabla_{\theta} \int_{\tau} P(\tau | \theta) R(\tau) && \text{Expand expectation} \\
 &= \int_{\tau} \nabla_{\theta} P(\tau | \theta) R(\tau) && \text{Bring gradient under integral} \\
 &= \int_{\tau} P(\tau | \theta) \nabla_{\theta} \log P(\tau | \theta) R(\tau) && \text{Log-derivative trick} \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P(\tau | \theta) R(\tau)] && \text{Return to expectation form} \\
 \therefore \nabla_{\theta} J(\pi_{\theta}) &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right] && \text{Expression for grad-log-prob}
 \end{aligned}$$

This is an expectation, which means that we can estimate it with a sample mean. If we collect a set of trajectories $\mathcal{D} = \{\tau_i\}_{i=1,\dots,N}$ where each trajectory is obtained by letting the agent act in the environment using the policy π_{θ} , the policy gradient can be estimated with

$$\hat{g} = \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau),$$

where $|\mathcal{D}|$ is the number of trajectories in \mathcal{D} (here, N).

This last expression is the simplest version of the computable expression we desired. Assuming that we have represented our policy in a way which allows us to calculate $\nabla_{\theta} \log \pi_{\theta}(a | s)$, and if we are able to run the policy in the environment to collect the trajectory dataset, we can compute the policy gradient and take an update step.

Implementing the Simplest Policy Gradient

We give a short Tensorflow implementation of this simple version of the policy gradient algorithm in [spinup/examples/pg_math/1_simple_pg.py](#). (It can also be viewed [on github](#).) It is only 122 lines long, so we highly recommend reading through it in depth. While we won't go through the entirety of the code here, we'll highlight and explain a few important pieces.

1. Making the Policy Network.

```

25 # make core of policy network
26 obs_ph = tf.placeholder(shape=(None, obs_dim), dtype=tf.float32)
27 logits = mlp(obs_ph, sizes=hidden_sizes+[n_acts])
28
29 # make action selection op (outputs int actions, sampled from policy)
30 actions = tf.squeeze(tf.multinomial(logits=logits, num_samples=1), axis=1)

```

This block builds a feedforward neural network categorical policy. (See the [Stochastic Policies](#) section in Part 1 for a refresher.) The `logits` tensor can be used to construct log-probabilities and probabilities for actions, and the `actions` tensor samples actions based on the probabilities implied by `logits`.

2. Making the Loss Function.

```
32 # make Loss function whose gradient, for the right data, is policy gradient
33 weights_ph = tf.placeholder(shape=(None,), dtype=tf.float32)
34 act_ph = tf.placeholder(shape=(None,), dtype=tf.int32)
35 action_masks = tf.one_hot(act_ph, n_acts)
36 log_probs = tf.reduce_sum(action_masks * tf.nn.log_softmax(logits), axis=1)
37 loss = -tf.reduce_mean(weights_ph * log_probs)
```

In this block, we build a “loss” function for the policy gradient algorithm. When the right data is plugged in, the gradient of this loss is equal to the policy gradient. The right data means a set of (state, action, weight) tuples collected while acting according to the current policy, where the weight for a state-action pair is the return from the episode to which it belongs. (Although as we will show in later subsections, there are other values you can plug in for the weight which also work correctly.)

⚠ You Should Know

Even though we describe this as a loss function, it is **not** a loss function in the typical sense from supervised learning. There are two main differences from standard loss functions.

1. The data distribution depends on the parameters. A loss function is usually defined on a fixed data distribution which is independent of the parameters we aim to optimize. Not so here, where the data must be sampled on the most recent policy.

2. It doesn’t measure performance. A loss function usually evaluates the performance metric that we care about. Here, we care about expected return, $J(\pi_\theta)$, but our “loss” function does not approximate this at all, even in expectation. This “loss” function is only useful to us because, when evaluated at the current parameters, with data generated by the current parameters, it has the negative gradient of performance.

But after that first step of gradient descent, there is no more connection to performance. This means that minimizing this “loss” function, for a given batch of data, has no guarantee whatsoever of improving expected return. You can send this loss to $-\infty$ and policy performance could crater; in fact, it usually will. Sometimes a deep RL researcher might describe this outcome as the policy “overfitting” to a batch of data. This is descriptive, but should not be taken literally because it does not refer to generalization error.

We raise this point because it is common for ML practitioners to interpret a loss function as a useful signal during training—"if the loss goes down, all is well." In policy gradients, this intuition is wrong, and you should only care about average return. The loss function means nothing.

⚠ You Should Know

The approach used here to make the `log_probs` tensor—creating an action mask, and using it to select out particular log probabilities—only works for categorical policies. It does not work in general.

3. Running One Epoch of Training.

```
45      # for training policy
46  def train_one_epoch():
47      # make some empty lists for logging.
48      batch_obs = []          # for observations
49      batch_acts = []         # for actions
50      batch_weights = []     # for R(tau) weighting in policy gradient
51      batch_rets = []        # for measuring episode returns
52      batch_lens = []        # for measuring episode lengths
53
54      # reset episode-specific variables
55      obs = env.reset()       # first obs comes from starting distribution
56      done = False            # signal from environment that episode is over
57      ep_rews = []             # list for rewards accrued throughout ep
58
59      # render first episode of each epoch
60      finished_rendering_this_epoch = False
61
62      # collect experience by acting in the environment with current policy
63      while True:
64
65          # rendering
66          if not(finished_rendering_this_epoch):
67              env.render()
68
69          # save obs
70          batch_obs.append(obs.copy())
71
72          # act in the environment
73          act = sess.run(actions, {obs_ph: obs.reshape(1,-1)})[0]
74          obs, rew, done, _ = env.step(act)
75
76          # save action, reward
77          batch_acts.append(act)
78          ep_rews.append(rew)
79
80          if done:
81              # if episode is over, record info about episode
82              ep_ret, ep_len = sum(ep_rews), len(ep_rews)
83              batch_rets.append(ep_ret)
84              batch_lens.append(ep_len)
85
86              # the weight for each logprob(a/s) is R(tau)
87              batch_weights += [ep_ret] * ep_len
88
89          # reset episode-specific variables
```

```

90         obs, done, ep_rews = env.reset(), False, []
91
92         # won't render again this epoch
93         finished_rendering_this_epoch = True
94
95         # end experience loop if we have enough of it
96         if len(batch_obs) > batch_size:
97             break
98
99         # take a single policy gradient update step
100        batch_loss, _ = sess.run([loss, train_op],
101                                feed_dict={
102                                    obs_ph: np.array(batch_obs),
103                                    act_ph: np.array(batch_acts),
104                                    weights_ph: np.array(batch_weights)
105                                })
106    return batch_loss, batch_rets, batch_lens

```

The `train_one_epoch()` function runs one “epoch” of policy gradient, which we define to be

1. the experience collection step (L62-97), where the agent acts for some number of episodes in the environment using the most recent policy, followed by
2. a single policy gradient update step (L99-105).

The main loop of the algorithm just repeatedly calls `train_one_epoch()`.

Expected Grad-Log-Prob Lemma

In this subsection, we will derive an intermediate result which is extensively used throughout the theory of policy gradients. We will call it the Expected Grad-Log-Prob (EGLP) lemma. [1]

EGLP Lemma. Suppose that P_θ is a parameterized probability distribution over a random variable, x . Then:

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0.$$

Proof

Recall that all probability distributions are *normalized*:

$$\int_x P_\theta(x) = 1.$$

Take the gradient of both sides of the normalization condition:

$$\nabla_\theta \int_x P_\theta(x) = \nabla_\theta 1 = 0.$$

Use the log derivative trick to get:

$$\begin{aligned}
0 &= \nabla_{\theta} \int_x P_{\theta}(x) \\
&= \int_x \nabla_{\theta} P_{\theta}(x) \\
&= \int_x P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x) \\
\therefore 0 &= \mathbb{E}_{x \sim P_{\theta}} [\nabla_{\theta} \log P_{\theta}(x)].
\end{aligned}$$

- [1] The author of this article is not aware of this lemma being given a standard name anywhere in the literature. But given how often it comes up, it seems pretty worthwhile to give it some kind of name for ease of reference.

Don't Let the Past Distract You

Examine our most recent expression for the policy gradient:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right].$$

Taking a step with this gradient pushes up the log-probabilities of each action in proportion to $R(\tau)$, the sum of *all rewards ever obtained*. But this doesn't make much sense.

Agents should really only reinforce actions on the basis of their *consequences*. Rewards obtained before taking an action have no bearing on how good that action was: only rewards that come *after*.

It turns out that this intuition shows up in the math, and we can show that the policy gradient can also be expressed by

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) \right].$$

In this form, actions are only reinforced based on rewards obtained after they are taken.

We'll call this form the "reward-to-go policy gradient," because the sum of rewards after a point in a trajectory,

$$\hat{R}_t \doteq \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}),$$

is called the **reward-to-go** from that point, and this policy gradient expression depends on the reward-to-go from state-action pairs.

⚠ You Should Know

But how is this better? A key problem with policy gradients is how many sample trajectories are needed to get a low-variance sample estimate for them. The formula we started with included terms for reinforcing actions proportional to past rewards, all of which had zero mean, but nonzero variance: as a result, they would just add noise to sample estimates of the policy gradient. By removing them, we reduce the number of sample trajectories needed.

An (optional) proof of this claim can be found [here](#), and it ultimately depends on the EGLP lemma.

Implementing Reward-to-Go Policy Gradient

We give a short Tensorflow implementation of the reward-to-go policy gradient in

[spinup/examples/pg_math/2_rtg_pg.py](#). (It can also be viewed [on github](#).)

The only thing that has changed from [1_simple_pg.py](#) is that we now use different weights in the loss function. The code modification is very slight: we add a new function, and change two other lines. The new function is:

```
12 def reward_to_go(rews):
13     n = len(rews)
14     rtgs = np.zeros_like(rews)
15     for i in reversed(range(n)):
16         rtgs[i] = rews[i] + (rtgs[i+1] if i+1 < n else 0)
17     return rtgs
```

And then we tweak the old L86-87 from:

```
86 # the weight for each Logprob(a/s) is R(tau)
87 batch_weights += [ep_ret] * ep_len
```

to:

```
93 # the weight for each Logprob(a_t/s_t) is reward-to-go from t
94 batch_weights += list(reward_to_go(ep_rews))
```

Baselines in Policy Gradients

An immediate consequence of the EGLP lemma is that for any function b which only depends on state,

$$\mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0.$$

This allows us to add or subtract any number of terms like this from our expression for the policy gradient, without changing it in expectation:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t) \right) \right].$$

Any function b used in this way is called a **baseline**.

The most common choice of baseline is the [on-policy value function](#) $V^\pi(s_t)$. Recall that this is the average return an agent gets if it starts in state s_t and then acts according to policy π for the rest of its life.

Empirically, the choice $b(s_t) = V^\pi(s_t)$ has the desirable effect of reducing variance in the sample estimate for the policy gradient. This results in faster and more stable policy learning. It is also appealing from a conceptual angle: it encodes the intuition that if an agent gets what it expected, it should “feel” neutral about it.

You Should Know

In practice, $V^\pi(s_t)$ cannot be computed exactly, so it has to be approximated. This is usually done with a neural network, $V_\phi(s_t)$, which is updated concurrently with the policy (so that the value network always approximates the value function of the most recent policy).

The simplest method for learning V_ϕ , used in most implementations of policy optimization algorithms (including VPG, TRPO, PPO, and A2C), is to minimize a mean-squared-error objective:

$$\phi_k = \arg \min_{\phi} \mathbb{E}_{s_t, \hat{R}_t \sim \pi_k} \left[\left(V_\phi(s_t) - \hat{R}_t \right)^2 \right],$$

where π_k is the policy at epoch k . This is done with one or more steps of gradient descent, starting from the previous value parameters ϕ_{k-1} .

Other Forms of the Policy Gradient

What we have seen so far is that the policy gradient has the general form

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \Phi_t \right],$$

where Φ_t could be any of

$$\Phi_t = R(\tau),$$

or

$$\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}),$$

or

$$\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t).$$

All of these choices lead to the same expected value for the policy gradient, despite having different variances. It turns out that there are two more valid choices of weights Φ_t which are important to know.

1. On-Policy Action-Value Function.

The choice

$$\Phi_t = Q^{\pi_\theta}(s_t, a_t)$$

is also valid. See [this page](#) for an (optional) proof of this claim.

2. The Advantage Function. Recall that the [advantage of an action](#), defined by $A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$, describes how much better or worse it is than other actions on average (relative to the current policy). This choice,

$$\Phi_t = A^{\pi_\theta}(s_t, a_t)$$

is also valid. The proof is that it's equivalent to using $\Phi_t = Q^{\pi_\theta}(s_t, a_t)$ and then using a value function baseline, which we are always free to do.

💡 You Should Know

The formulation of policy gradients with advantage functions is extremely common, and there are many different ways of estimating the advantage function used by different algorithms.

💡 You Should Know

For a more detailed treatment of this topic, you should read the paper on [Generalized Advantage Estimation](#) (GAE), which goes into depth about different choices of Φ_t in the background sections.

That paper then goes on to describe GAE, a method for approximating the advantage function in policy optimization algorithms which enjoys widespread use. For instance, Spinning Up's implementations of VPG, TRPO, and PPO make use of it. As a result, we strongly advise you to study it.

Recap

In this chapter, we described the basic theory of policy gradient methods and connected some of the early results to code examples. The interested student should continue from here by studying how the later results (value function baselines and the advantage formulation of policy gradients) translate into Spinning Up's implementation of [Vanilla Policy Gradient](#).