

On the policy gradient method

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license.

Kamila Zdybał

kamilazdybal.github.io, kamila.zdybal@gmail.com

Please feel free to contact me with any suggestions, corrections or comments.

Preface

Reinforcement learning (RL) is all about finding optimal *policies*—rules of acting on or behaving in an environment. The policy, π , is simply a function that takes observations of the environment as inputs and it outputs actions to be executed in that environment. There’s a couple of different methodological approaches to how one can obtain the optimal policy (or *train* the RL). In particular, the policy gradient method (PGM) is one of the RL algorithms where we directly learn the policy function, π_{θ} , with its trainable parameters $\theta \in \mathbb{R}^d$. The functional form for the policy can be imposed by us in any way that we want, as long as it is differentiable¹ with respect to θ . Since the advent of deep learning, we commonly use artificial neural networks as the function approximators for π_{θ} . In PGM, π is a very specific function and one that outputs **the probability of each action being selected**². To actually select an action, we sample from those probabilities. In other words, we write that $\pi = \pi(a_t | s_t)$. With sufficient training, the optimal PGM policy returns the maximum probability for the *best* action to be taken in each state in the environment.

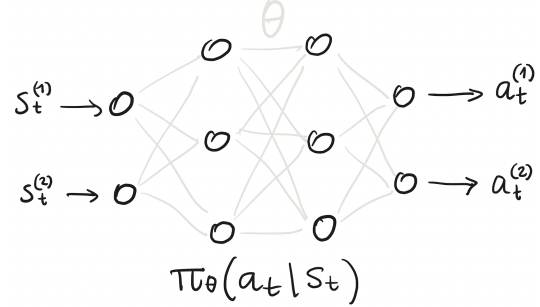
This document is an educational tutorial deriving and explaining the PGM in more depth. The pre-requisites are that you understand what RL is in general (what is state, action, reward) and that you’ve been exposed to training of artificial neural networks (ANNs). I specifically wanted to root the intuition behind PGM in the intuition that you may already have from training general-purpose ANNs. This document is a bit mathy, but stick with it and you will find that it’s actually not that much more of a conceptual advance from the math of gradient descent!

Keywords

reinforcement learning, policy gradient method, deep neural networks

¹The differentiability requirement is necessary because the core mechanism of the PGM involves computing the gradient of π_{θ} , as we shall see later. Hence the algorithm is called the *policy gradient* method.

²This is an important assumption that we make use of in the derivation of the PGM. If the policy ever outputs anything else, the present derivation of the PGM may need a modification.



The core concept

The PGM, at its core, is trained in the same way as we train artificial neural networks in general. We compute the gradient of some performance measure with respect to the trainable parameters, θ , and we update θ in the direction dictated by that gradient by taking a small step (measured by the learning rate, α) in that direction.

In a classic regression problem, the performance measure could be the mean-squared-error between the true and the predicted targets, call it $\text{MSE}(\theta)$. We want to nudge θ in the direction that *minimizes* $\text{MSE}(\theta)$, so we would be taking a step in the direction *opposite* to the gradient of $\text{MSE}(\theta)$. The parameter update rule in that case is

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta} \text{MSE}(\theta). \quad (1)$$

Doing this update iteratively gives us the classic gradient *descent* algorithm.

To be completely precise, in the deep learning world, we use a *stochastic* gradient descent, which means that $\text{MSE}(\theta)$ is estimated from samples of true vs. predicted targets. We cannot know the raw $\text{MSE}(\theta)$ for *all* possible values of targets—there’s astronomically many of them, if not infinitely many. We can only estimate $\text{MSE}(\theta)$ from those values of targets that we have actually sampled as our “training data”. But that’s okay. As long as our data sampling was statistically representative, the MSE that we compute from those samples is a good enough approximation of the MSE that would have been computed on all possible samples—this is the principle behind the Monte Carlo estimates. We denote this stochastic estimate of $\text{MSE}(\theta)$ as $\widehat{\text{MSE}}(\theta)$ and so, instead of Eq. (1), we actually compute

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta} \widehat{\text{MSE}}(\theta). \quad (2)$$

Conversely, in the PGM, our performance measure, call it $J(\theta)$, is constructed (in some smart way) from the rewards given to the RL agent in the environment. We want to maximize the rewards in the long run, right? So we should be taking steps

in the direction that *maximizes* $J(\theta)$, or in the direction of the gradient of $J(\theta)$. Hence, the parameter update rule in PGM is

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_{\theta} J(\theta). \quad (3)$$

As opposed to training the regressor, this is the gradient *ascent* algorithm.

So far so good. If we can compute, or even reasonably approximate $\nabla_{\theta} J(\theta)$, then we are all set. But the name of the game in PGMs is how do we get $\nabla_{\theta} J(\theta)$ in practice? In the next sections, I would like to first show you why getting the raw $\nabla_{\theta} J(\theta)$ is impossible in practice (for much the same reasons why we couldn't get the raw $\text{MSE}(\theta)$) and then show you a bunch of neat math rearrangements that make estimating $\nabla_{\theta} J(\theta)$ possible!

How do we construct the performance measure, $J(\theta)$?

Before we construct a specific $J(\theta)$, let's take a step back and see what options and restrictions we are dealing with.

- #1 We've already said that it makes sense that $J(\theta)$ is *somehow* constructed from the rewards given to the RL agent as he navigates the environment. In the end, the essence of RL is learning from experience by receiving feedback signals (rewards), which are the only mechanism through which we correct and improve the behavior of the RL agent in the future. We expect that the optimal policy, π_{θ}^* , leads to solving the task in the environment perfectly, which we can recognize by observing that the agent receives the highest rewards possible. Usually in RL we use the discounted future rewards, but there is some leeway to how exactly we can make the rewards enter the performance measure.
- #2 We have to be able to estimate $\nabla_{\theta} J(\theta)$ from actual rollouts in the environment. There's no other way. We have to set the agent running and interacting with the environment and see what he accomplishes. We do not have any other mechanism for obtaining the "training data". In that sense, similarly as in all neural network training, we will be computing stochastic estimates of $\nabla_{\theta} J(\theta)$, which we denote as $\widehat{\nabla_{\theta} J(\theta)}$. So, instead of Eq. (3), in PGM we actually compute

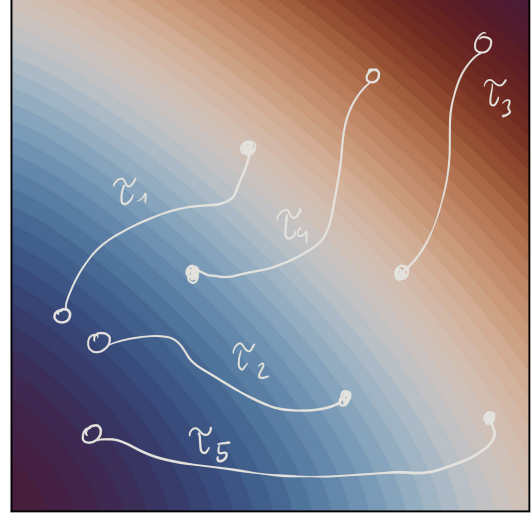
$$\theta_{t+1} \leftarrow \theta_t + \alpha \widehat{\nabla_{\theta} J(\theta)}, \quad (4)$$

for reasons analogous to those that made us arrive at the stochastic estimate of the $\text{MSE}(\theta)$ ³.

Okay, so let's use the reasoning from #1 and write out what $J(\theta)$ could be. The most mathematically general statement would be to say that

$$J(\theta) := \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)]. \quad (5)$$

³For the moment, it is not clear why we are after the stochastic estimate of the **gradient** of $J(\theta)$ and not the stochastic estimate of $J(\theta)$ directly. This is where Eq. (4) differs from Eq. (2). Please be patient, as this will become clear later!



Here, τ is the trajectory that the RL agent takes in the environment, sampled using the current version of the policy, π_{θ} , when its parameters are θ . $R(\tau)$ is the "total return" from following the trajectory τ and it is *somehow* constructed from the rewards in the environment, perhaps like so:

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t, \quad (6)$$

where $0 < \gamma < 1$ is the discount factor and r_t is the reward value given to the agent at the time step t . The sum then takes account of all incremental rewards that were obtained throughout the duration of the trajectory τ , i.e., from time step $t = 0$ to some final time step $t = T$.

Let's ponder about that expected value, $\mathbb{E}_{\tau \sim \pi_{\theta}}$, for a bit. First, we could also expand the definition of the expected value (in the discrete setting) and write Eq. (5) differently as

$$J(\theta) = \sum_{i=0}^N [R(\tau_i) \cdot P(\tau_i | \theta)], \quad (7)$$

In other words, the expected value across all sampled trajectories is the sum of total returns weighted by the probability, P , of the specific trajectory τ_i occurring under the current policy parameters θ . N can be a *really* huge number. In fact, it can be infinite and we would then write the expected value in the continuous setting as

$$J(\theta) = \int_{\tau} [R(\tau) \cdot P(\tau | \theta)] d\tau. \quad (8)$$

But for the moment, the expected value itself is not scary. We could, in the end, estimate it from samples of a couple of rollouts actually taken. We could let the agent take a few, or even many trajectories, τ , and we could average the total returns across those to give us some estimate of $J(\theta)$. This would give us a Monte Carlo estimate of $J(\theta)$,

$$J(\theta) \approx \widehat{J(\theta)} = \frac{1}{N} \sum_{i=1}^N R(\tau_i), \quad (9)$$

and the higher N the better this estimate would be. In other words, getting $\widehat{J}(\theta)$ wouldn't be difficult at all if we ever needed that. But recall that in Eq. (3) we need the gradient of $J(\theta)$ in order to know how to improve the policy. And things complicate when we take that gradient...

The complication that arises when taking the gradient of $J(\theta)$

Let's take the gradient of $J(\theta)$ as defined in Eq. (4) to see what happens! First, since the gradient is with respect to the parameters θ , we ask if the expected value is a function of θ ? Well, yes! It is, because the probability of a given trajectory occurring depends on the current policy (which guides the dynamics of the agent's movement across the environment). And the policy explicitly depends on θ . Note that $R(\tau)$ is not a function of the policy nor its parameters—the way that rewards are placed in the environment is simply a function of the environment itself. Eq. (7) shows explicitly that there is a dependence on θ , namely, the probability P is a function of θ . So, at most we can write that

$$\nabla_{\theta} J(\theta) = \sum_{i=0}^N [R(\tau_i) \cdot \nabla_{\theta} P(\tau_i | \theta)] . \quad (10)$$

Computing (or estimating) the right-hand-side of the above equation is impossible in practice because we have no way of knowing how the probability P changes as we vary the parameters θ by just sampling a few trajectories. We'd have to sample *all of them* to know $\nabla_{\theta} P(\tau | \theta)$ because those probabilities are coupled. In other words, we cannot compute how the probability of one specific trajectory arising changes with θ in isolation, because it depends on how the probabilities of all the other trajectories shift in response. We would need to understand the entire distribution's response to θ . The current form of the gradient of $J(\theta)$ violates our restriction #2. But a beautiful simplification arises when we use one math trick and also expand a bit what that probability P is equal to!

And while we're on it, note that we also cannot obtain the gradient of $\widehat{J}(\theta)$ like we do with the $MSE(\theta)$ in classic regression problems. If $\widehat{J}(\theta)$ is obtained with a Monte Carlo estimate, like in Eq. (9), then it only depends on the total returns, $R(\tau_i)$, which aren't explicitly functions of θ —they depend on the policy and its parameters only through random sampling. We couldn't differentiate through the random sample. So our last hope is going after $\nabla_{\theta} \widehat{J}(\theta)$.

Where maths comes to the rescue

First, we are going to use the “log-derivative trick”. It states that

$$\nabla_{\theta} P(\tau | \theta) = P(\tau | \theta) \cdot \nabla_{\theta} \ln(P(\tau | \theta)) . \quad (11)$$

This simply comes from rearrangement⁴ of what the derivative of a natural logarithm of a multivariate function is equal to. We could plug in the result from Eq. (11) into Eq. (10) but that doesn't improve our situation, we still ultimately have to deal with the gradient of that probability P (or the natural logarithm of it). Let's see if we can do something about that P .

What is $P(\tau | \theta)$, really?

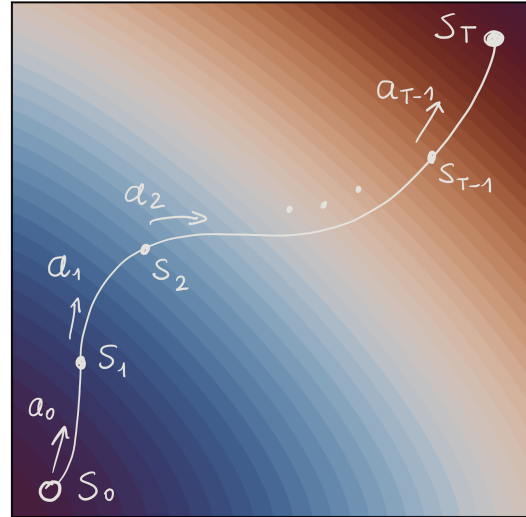
How would we compute $P(\tau | \theta)$ for one specific sampled trajectory, τ , given that τ is the following sequence of states, s_t , and actions, a_t ,

$$\tau = \{s_0, a_0, s_1, a_1, s_2, a_2, \dots, s_{T-1}, a_{T-1}, s_T\} . \quad (12)$$

Well, three factors contribute to *this particular* trajectory (sequence of states and actions) arising and we list them in the form of three questions:

1. What was the chance that τ started in *this particular* initial state, s_0 ?
2. What was the chance that the action selected in s_0 was *this particular* a_0 ?
3. What was the chance that by executing this particular a_0 in this particular s_0 we then entered *this particular* next state, s_1 ?

and then, you can continue asking yourself the questions 2. and 3. in a loop as you traverse all the states in τ , all the way till you get to the terminal state, s_T !



So let's write these three factors mathematically in terms of probabilities:

⁴For an independent variable, x , we have that $\frac{d}{dx} \ln(x) = \frac{1}{x}$, whereas for a multivariate function f , the chain rule gives $\nabla \ln(f) = \frac{1}{f} \cdot \nabla f$. Now just multiply the latter by f and you get the “log-derivative trick”.

1. The probability that τ started in *this particular* initial state, s_0 , is equal to $P_0(s_0)$, where P_0 is some distribution over the possible starting states in the environment. For instance, for entirely fair environments, like a 2D grid world, all grid cells may be equally likely, so $P_0(s_0) \sim \mathcal{U}$. But some environments, like games, may have more likely starting positions, and some starting positions are impossible. You may even implement a completely deterministic starting position so that all τ originate from the same location. Sky's the limit! But the important take-away is that $P_0(s_0)$ is the property of the environment, and it is not a function of the policy!
2. The probability that the action selected in s_0 was *this particular* a_0 is simply equal to $\pi_{\theta}(a_0|s_0)$ —our policy explicitly outputs that probability! (Recall what I mentioned in the preface about π needing to output probabilities explicitly. This is where we are using this assumption!)
3. The probability that by executing this particular a_0 in this particular s_0 we then entered *this particular* next state, s_1 , is equal to $P_t(s_1|s_0, a_0)$, where P_t is also known as the state transition probability. This probability is the property of the environment, too. More specifically, a property of any *stochastic* environment. A deterministic environment would have each $P_t = 1$. It defines the dynamics of the environment. It is also not a function of the policy!

With all of this, we are ready to write what is $P(\tau|\theta)$ equal to:

$$P(\tau|\theta) = P_0(s_0) \cdot \prod_{t=0}^{T-1} \pi_{\theta}(a_t|s_t) \cdot P_t(s_{t+1}|s_t, a_t). \quad (13)$$

Note that we simply multiply all of those probabilities because all of those events are independent from one another.

Putting it all together

The mathematical beauty comes from substituting Eq. (13) into the gradient-of-the-logarithm term in Eq. (11)!

First, if we take the natural logarithm of Eq. (13), all products become sums of logarithms:

$$\begin{aligned} \ln(P(\tau|\theta)) &= \ln(P_0(s_0)) + \\ &+ \sum_{t=0}^{T-1} \ln(\pi_{\theta}(a_t|s_t)) + \\ &+ \sum_{t=0}^{T-1} \ln(P_t(s_{t+1}|s_t, a_t)). \end{aligned} \quad (14)$$

Second, if we now take the gradient of Eq. (15), the first and the third term on the right-hand-side vanish, because as we've reasoned earlier, they are not a function of the policy parameters, θ ! So we are left with

$$\nabla_{\theta} \ln(P(\tau|\theta)) = \sum_{t=0}^{T-1} \nabla_{\theta} \ln(\pi_{\theta}(a_t|s_t)). \quad (15)$$

This is a **really key result**! It allows us to shift from the need of computing gradients of those probabilities of trajectories,

P , which we've mentioned are intractable, to only needing to compute gradients of the policy function itself! (The natural logarithm of it, that is.) And that we can do because we know the functional form of the policy. If it is a neural network, we can always compute how action probabilities change with a specific change in θ . In fact, this is what automatic differentiation will keep track of anyway for us when we use a package like PyTorch or TensorFlow.

Alright, but let's write out the complete version of $\nabla_{\theta} J(\theta)$ to see what we have accomplished. We have obtained

$$\nabla_{\theta} J(\theta) = \sum_{i=0}^N [R(\tau_i) \cdot P(\tau_i|\theta) \cdot \nabla_{\theta} \ln(P(\tau_i|\theta))] \quad (16)$$

from just using the “log-derivative trick”. And now we also have

$$\nabla_{\theta} J(\theta) = \sum_{i=0}^N \left[R(\tau_i) \cdot P(\tau_i|\theta) \cdot \sum_{t=0}^{T-1} \nabla_{\theta} \ln(\pi_{\theta}(a_t|s_t)) \right]. \quad (17)$$

Note that we still have the probability P appearing in this equation, but not its gradient. In fact, due to this probability appearing explicitly, we can again recognize the definition of an expected value in Eq. (17)! Hence, we can write that

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[R(\tau) \cdot \sum_{t=0}^{T-1} \nabla_{\theta} \ln(\pi_{\theta}(a_t|s_t)) \right]. \quad (18)$$

The result in Eq. (18) allows $\widehat{\nabla_{\theta} J(\theta)}$ to be entirely computable from just samples of selected trajectories. This is achievable with a Monte Carlo estimate:

$$\widehat{\nabla_{\theta} J(\theta)} \approx \widehat{\nabla_{\theta} J(\theta)} = \frac{1}{N} \sum_{i=1}^N \left[R(\tau_i) \cdot \sum_{t=0}^{T-1} \nabla_{\theta} \ln(\pi_{\theta}(a_t|s_t)) \right]. \quad (19)$$

In other words, we can estimate this present expected value by running a couple of rollouts in the environment and averaging them. For each rollout, at each time step, we will also be computing the gradient of the policy term.

Take a moment to note that we couldn't compute the raw $\nabla_{\theta} J(\theta)$, we also couldn't compute $\nabla_{\theta} \widehat{J(\theta)}$, but we now know how to compute $\widehat{\nabla_{\theta} J(\theta)}$ and that's the success story of the PGM!

Coding the policy gradient method in PyTorch

I will now show you the simplest possible example of setting up and training the PGM using PyTorch, given everything that we have derived thus far. This will help ground the earlier reasoning and derivations, and you will see that the code is pretty much a verbose translation of the maths we've been discussing!

Here's a few imports that we'll need:

```
import numpy as np
import torch
import torch.nn as nn
from torch.distributions import Categorical
```

The PGM code is largely independent of the environment that you have⁵, but we assume that your environment class has functions `reset()` and `step()` with specific returned values and that there is some trajectory termination criterion such that we can step in the environment in a `while` loop and it won't loop forever. You can initialize your custom environment:

```
env = CustomEnv()
```

Let's start by constructing a deep neural network that will serve as our function approximator to π_θ . First, you can use any architecture that you'd like! In the module defined below, we use two hidden layers with eight neurons each and hyperbolic tangent activations. The input state has only one element and we allow two actions.

```
class Policy(nn.Module):
    def __init__(self):
        super().__init__()

        self.net = nn.Sequential(
            nn.Linear(1, 8),
            nn.Tanh(),
            nn.Linear(8, 8),
            nn.Tanh(),
            nn.Linear(8, 2))

    def forward(self, x):
        logits = self.net(x)

        return Categorical(logits=logits)
```

I'd like you to pay special attention that this network **needs to output probabilities** to satisfy our primary assumption about π_θ . Remember that if $\pi_\theta(a_t|s_t)$ didn't return probability directly, the very derivation of the PGM would not have worked out. One way to implement this is to wrap the outputs of the

last neural layer with what we call a "softmax" transformation. The softmax is a really neat function that takes a vector of a bunch of real numbers, called "logits", and turns them into probabilities according to their magnitude, and in such a way that all of those probabilities always sum to unity. Let's say that the outputs of the last layer in our policy neural network are denoted h_i (where each h_i is a function of the current policy parameters, θ , and the current input state, s_t) and we have n of them. The softmax transformation of each h_i is

$$\text{softmax}(h_i(s_t, \theta)) = \frac{e^{h_i(s_t, \theta)}}{\sum_{j=1}^n e^{h_j(s_t, \theta)}}, \quad (20)$$

with the property that

$$\sum_{i=1}^n \text{softmax}(h_i(s_t, \theta)) = 1. \quad (21)$$

If we make the policy return that softmax for each action i ,

$$\pi_\theta(a_t^{(i)}|s_t) = \text{softmax}(h_i(s_t, \theta)), \quad (22)$$

then we're all set! In practice, when our actions are discrete (such as go left or go right), we can make the policy network return a categorical distribution over what the last neural layer computes, like we've done here. The categorical distribution from PyTorch runs the softmax over the logits under the hood when we sample actions from that distribution.

We can already initialize the policy network:

```
policy = Policy()
```

In order to compute the Monte Carlo estimates as per Eq. (19), we write a function that generates one trajectory, τ . Later, we can run it a couple of times inside the training loop to improve our Monte Carlo estimate. This function steps in the environment, each time collecting incremental rewards, r_t , and computing the natural log of the policy output for the action selected, $\ln(\pi_\theta(a_t|s_t))$. After the trajectory has terminated, we compute the total return as per Eq. (6), given the user-specified discount factor, γ .

```
def generate_τ(env,
              policy,
              γ=0.95):
    # Reset the environment to an initial state, s_0:
    state = env.reset()

    ln_π_list = []
    rewards = []
    done = False

    # Step in the environment until termination:
    while not done:
        # Get the output of π for the current state:
        distribution = policy(state.unsqueeze(0))
```

⁵The main assumption about the environment is that there is one input state value and two output actions in the policy function. You can always change the policy network architecture to accommodate a different number of state elements and actions.

```

# Sample action from the current  $\pi$ :
action = distribution.sample()

# Compute the  $\ln$  of  $\pi$  for that action:
ln_pi = distribution.log_prob(action)

# Take a step in the environment:
state, r, done = env.step(action.item())

# Save the current  $\ln(\pi)$  and reward:
ln_pi_list.append(ln_pi)
rewards.append(r)

# Once  $\tau$  has terminated: - - - - -

# Compute the total return from this  $\tau$ :
R = sum([(gamma ** t) * r for t, r in
        enumerate(rewards)])
R = torch.tensor(R, dtype=torch.float32)

# Compute the sum of  $\ln(\pi)$  across the trajectory:
sum_ln_pi = torch.stack(ln_pi_list).sum()

return sum_ln_pi, R

```

Now we're ready to code the training loop! Let's start with defining a few numbers:

```

n_episodes = 1000
N = 10
alpha = 1e-3
gamma = 0.95

```

We use the Adam gradient descent optimizer, where the parameters to optimize are the policy parameters, θ :

```

theta = policy.parameters()
optimizer = torch.optim.Adam(theta, lr=alpha)

```

And here's the full training loop:

```

for _ in range(1, n_episodes+1):
    optimizer.zero_grad()

    loss = 0.0

    # Monte Carlo estimate from N trajectories:
    for _ in range(0, N):
        sum_ln_pi, R = generate_tau(env, policy, gamma)

        loss += -(R.detach() * sum_ln_pi)

    # Compute the arithmetic average:
    loss /= N

    # Backpropagate loss and update  $\theta$ :
    loss.backward()
    optimizer.step()

```

The great thing about using PyTorch's automatic differentiation is that all that we need to compute ourselves is this *loss function* (notice no gradient yet on the $\ln(\pi_{\theta}(a_t|s_t))$ term):

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \left[R(\tau_i) \cdot \sum_{t=0}^{T-1} \ln(\pi_{\theta}(a_t|s_t)) \right]. \quad (23)$$

When we backpropagate this loss, *i.e.*, run `loss.backward()`, PyTorch will be computing the gradient of it internally for us, and hence we will get

$$\nabla_{\theta} \mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \left[R(\tau_i) \cdot \sum_{t=0}^{T-1} \nabla_{\theta} \ln(\pi_{\theta}(a_t|s_t)) \right] \quad (24)$$

for free, which is the same as Eq. (19) but with flipped sign so that we enforce gradient *ascent* instead of gradient *descent*. In other words, we never have to code the computation of the gradient ourselves, isn't that great? But there's one caveat, we need to make sure that $R(\tau_i)$ is always a constant and it does not take part in gradient computation—we only want the gradient of the $\ln(\pi_{\theta}(a_t|s_t))$ term, right? That's why we run `R.detach()` to detach this reward tensor from the computational graph.

And that's pretty much it! All that will likely remains for the user is to tune the hyper-parameters while monitoring the training outcomes.

On exploration in PGM

I wanted to make a comment on how the idea of exploration vs. exploitation in PGM is baked in the concept of using a probabilistic policy, π . In other RL algorithms, such as in deep Q-learning, the user explicitly specifies the exploration probability at each rollout in the environment. Actions are either taken at random during exploration or are taken according to the current policy during exploitation. But here we didn't use the exploration probability explicitly. Instead, the action sampling step, `action = distribution.sample()`, allows for sub-optimal actions to be selected from time to time, because even if policy improves, there is still the element of randomness coming from random sampling. In addition, as the policy improves with training, we would hope that the probabilities for the best actions increase significantly above the probabilities for all the remaining actions. That way, we naturally get the effect of "exploration probability decay" with training time, which in deep Q-learning the user has to implement explicitly.