

Week 5: Policy Optimization I: Policy Gradient and Actor-Critic

Bolei Zhou

UCLA

February 3, 2026

Announcement

- ① Assignment 2 due by the end of this week

This Week's Plan

- ① Introduction on policy function versus value function
- ② Policy-based reinforcement learning
- ③ Monte-Carlo policy gradient
- ④ Actor-Critic methods

Value-based RL vs Policy-based RL

- ① In previous lectures, we focused on value-based RL and had value function approximation with parameter θ

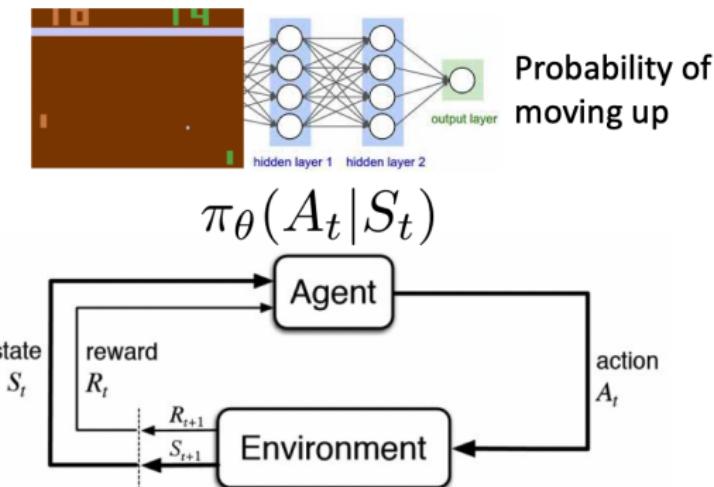
$$V_\theta(s) \approx v^\pi(s)$$
$$Q_\theta(s, a) \approx q^\pi(s, a)$$

- ② A policy is generated directly from the value function using ϵ -greedy or greedy $a_t = \arg \max_a Q(a, s_t)$
- ③ Instead, we can also parameterize the policy function as $\pi_\theta(a|s)$ where θ is the learnable policy parameter

RL Diagram in the View of Policy Optimization

no need for a state or q value function

- ① Action from the policy is all we need, then let's optimize the policy directly!



Value-based RL versus Policy-based RL

① Value-based RL

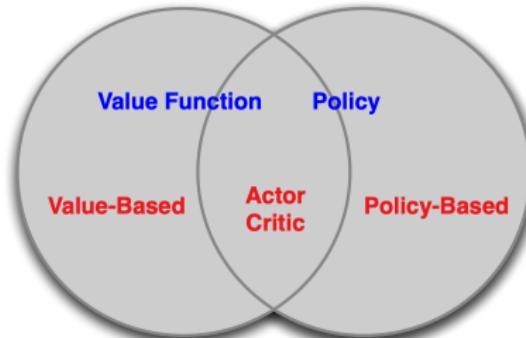
- ① learn value function
- ② implicit policy based on the value function

② Policy-based RL

- ① no value function
- ② learn policy directly

③ Actor-critic

- ① learn both policy and value function



Strengths of Policy-based RL

① Strengths:

- ① Better convergence properties: we are guaranteed to converge on a local optimum (worst case) or global optimum (best case)
- ② Policy gradient is more effective in a high-dimensional action space
- ③ Policy gradient can learn stochastic policies, while value function can't
probability over actions

② Weaknesses:

- ① typically converges to a local optimum
- ② Evaluating a policy is inefficient and has high variance

Two Types of Policies

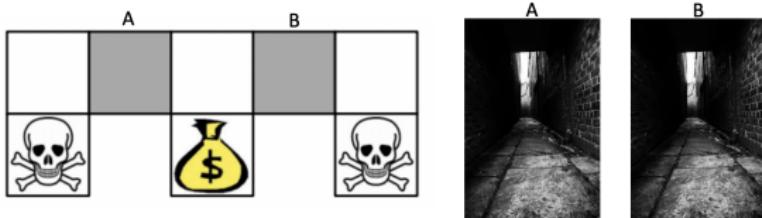
- ① Deterministic: given a state, the policy returns the action to take (discrete action space) or returns a continuous value (continuous action space)
- ② Stochastic: given a state, the policy returns a probability distribution of the actions (e.g., 40% chance to turn left, 60% chance to turn right), or a probability density function for continuous action

Example: Rock-Paper-Scissor



- ① Two-player game
- ② What is the best policy?
 - ① It is easy to beat a deterministic policy
 - ② A uniform random policy is optimal (Nash equilibrium)

Example: Alleyways in the Gridworld



- ① The agent cannot differentiate the two grey states (the two alleyways look the same to the agent)
- ② Considers the following features (for all N, E, S, W)

$$\psi(s, a) = [\mathbf{1}(\text{wall to N}, a = \text{move E}), \mathbf{1}(\text{wall to S}, a = \text{move W}), \dots]$$

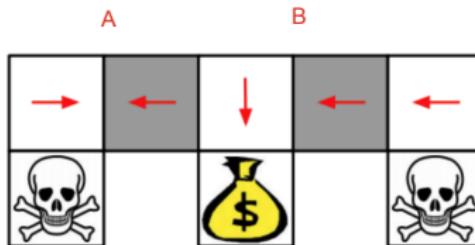
- ③ If it is value-based RL, value function approximation as:
same features will represent the two states

$$Q_\theta(s, a) = f(\psi(s, a), \theta)$$

- ④ If it is policy-based RL, policy function approximation as:

$$\pi_\theta(a|s) = g(\psi(s, a), \theta)$$

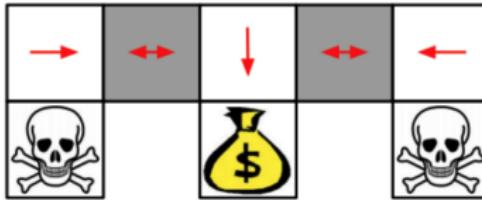
Example: Alleyways in the Gridworld



in value based, A and B must be the same policy (based on the value function?)

- ① Value-based RL learns a near-deterministic policy, e.g., greedy or ϵ -greedy
- ② Because of the aliasing (the agent cannot differentiate two states), an optimal deterministic policy from value-based RL will either
 - ① move W in both grey states (shown by red arrows)
 - ② move E in both grey states
- ③ Either way, 50% chance will get stuck

Example: Alleyways in the Gridworld



- ① Policy-based RL can learn the optimal stochastic policy
- ② An optimal stochastic policy will randomly move E or W in two grey states

$$\pi_{\theta}(\text{move E} | \text{wall to N and S}) = 0.5$$

$$\pi_{\theta}(\text{move W} | \text{wall to N and S}) = 0.5$$

- ③ For any starting point, it will reach the goal state in a few steps with high probability

Objective of Optimizing Policy

- ① Objective: Given a policy approximator $\pi_\theta(a|s)$ with parameter θ , find the best θ
want to optimize cumulative reward
- ② How do we measure the quality of a policy π_θ ?
plug in policy in the environment and do the roll out
- ③ In episodic environments we can use the start value

$$J_1(\theta) = V^{\pi_\theta}(s_1)$$

s1 = initial state

- ④ In continuing environments

- ① we can use the average value

probability of choosing the state given the policy
value of starting at s?

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

get expected state value over all states

- ② or the average reward per time-step

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) R(s, a)$$

where d^{π_θ} is stationary distribution of Markov chain for π_θ

Objective of Optimizing Policy

- ① The value of the policy is defined as

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_t r(s_t, a_t) \right]$$

use per step rewards to evaluate quality of a policy
sampling actions from our policy

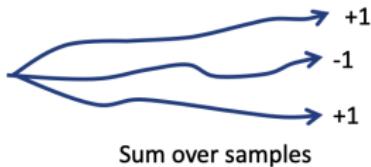
objective function measures the performance of whole trajectories sampled from our policy

run m times

$$\approx \frac{1}{m} \sum_m \sum_t r(s_{t,m}, a_{t,m})$$

m indexes each trajectory

It is the same as the value function we defined in the value-based RL



- ① τ is a trajectory sampled from the policy function π_θ
- ② The goal of policy-based RL

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_t r(s_t, a_t) \right]$$

tau = [a1, a2, a3, ...] (whole trajectories)

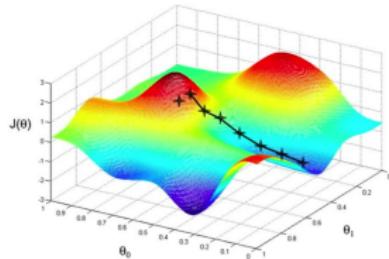
tau are the actions sampled from our policy

Objective of Optimizing Policy

- ① Policy-based RL is an optimization problem that finds θ that maximizes $J(\theta)$
- ② If $J(\theta)$ is differentiable, we can use gradient-based methods:
 - ① gradient ascend
 - ② conjugate gradient
 - ③ quasi-newton
- ③ If $J(\theta)$ is non-differentiable or hard to compute the derivative, some derivative-free black-box optimization methods:
 - ① Cross-entropy method (CEM)
 - ② Evolution strategy
 - ③ Finite difference

can't calculate the derivative $df(x)/dx$

Policy Optimization using Gradient Ascend



- ① Consider a function $J(\theta)$ to be any policy objective function
- ② Goal is to find parameter θ^* that maximizes $J(\theta)$ by ascending the gradient of the policy, w.r.t parameter θ

$$\Delta\theta = \alpha \nabla_{\theta} J(\theta)$$

- ③ Adjust θ in the direction of the gradient, where α is step-size
- ④ Define the gradient of $J(\theta)$ to be

$$\nabla_{\theta} J(\theta) = \left(\frac{\partial J(\theta)}{\partial \theta_1}, \frac{\partial J(\theta)}{\partial \theta_2}, \dots, \frac{\partial J(\theta)}{\partial \theta_n} \right)^T$$

Policy Optimization using Derivative-free Methods

- ① Sometimes we cannot compute the derivative, i.e., $\nabla_{\theta} J(\theta)$
- ② Derivative free methods:
 - ① Evolution Strategy (ES)
 - ② Cross Entropy Method (CEM)
 - ③ Finite Difference

Derivative-free Method: Evolution Strategy (ES)

① $\theta^* = \arg \max J(\theta)$

$\backslash\pi_\\theta$ is differentiable?

- ② Treat $J(\theta)$ as a black box score function (not differentiable)

this helps maximize $J(\theta)$ based on changing θ

Algorithm 1 Evolution Strategy

```
1: for iter  $t = 0, 1, 2, \dots$  do      generate perturbation (noise), add to policy parameters  
2:   sample  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, I)$  to generate perturbed  $\theta_i = \theta_t + \delta\epsilon_i$   
3:   execute roll-out under  $\theta_i$  for  $i = 1, \dots, n$  to compute  $J(\theta_i)$ play out  
4:   set  $\theta_{t+1} \leftarrow \theta_t + \alpha \frac{1}{n\delta} \sum_{i=1}^n J(\theta_i) \epsilon_i$           and get the sample return  
5: end for           use a weighted sum with noise to update parameters
```

- ③ ES is easy to implement and scale. Running on a computing cluster of 80 machines and 1,440 CPU cores, the implementation is able to train a 3D MuJoCo humanoid walker in only 10 minutes
 - ① <https://openai.com/blog/evolution-strategies/>
- ④ Example code: https://github.com/ucla-rlcourse/RLexample/blob/master/derivativefree/es_cem_agent.py

Derivative-free Method: Cross-Entropy Method

- ① $\theta^* = \arg \max J(\theta)$
- ② Treat $J(\theta)$ as a black box score function (not differentiable)

Algorithm 2 CEM

```
1: for iter  $t = 1, 2, 3, \dots$  do
2:    $\mathcal{C} = \{\}$ 
3:   for parameter set  $e = 1$  to  $n$  do                                think  $N(\mu, \sigma^2)$ 
4:     sample  $\theta^{(e)} \sim P_{\mu^{(t)}}(\theta)$  probability distribution over parameters
5:     execute roll-outs under  $\theta^{(e)}$  to evaluate  $J(\theta^{(e)})$ 
6:     store  $(\theta^e, J(\theta^{(e)}))$  in  $\mathcal{C}$            complete the roll outs
7:   end for                                         update mu to increase log likelihood of generating a good set
8:    $\mu^{(t+1)} = \arg \max_{\mu} \sum_{k \in \hat{\mathcal{C}}} \log P_{\mu}(\theta^{(k)})$           based on choosing only the theta params
   where  $\hat{\mathcal{C}}$  are the top 10% of  $\mathcal{C}$  ranked by  $J(\theta^{(e)})$  with the best performance
9: end for
```

-
- ③ Example code: [https://github.com/ucla-rlcourse/RLexample/
blob/master/derivativefree/es_cem_agent.py](https://github.com/ucla-rlcourse/RLexample/blob/master/derivativefree/es_cem_agent.py)

Approximate Gradients by Finite Difference

- ① To evaluate policy gradient of $\pi_\theta(a|s)$
- ② For each dimension $k \in [1, n]$
 - ① estimate k th partial derivative of objective function by perturbing θ by a small amount δ in k th dimension

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \delta u_k) - J(\theta)}{\delta}$$

where u_k is unit vector with 1 in k th component, 0 elsewhere

- ③ use n evaluations to compute policy gradient in total n dimensions
- ④ Though noisy and inefficient, it works for arbitrary policies, even if policy is not differentiable.

Compute the Policy Gradient Analytically

REVIEW THIS SLIDE GANG

- ① π_θ is a probability function or density function with parameter θ
- ② Assume policy π_θ is differentiable whenever it is no-zero
- ③ We can compute the gradient $\nabla_\theta \pi_\theta(a|s)$
- ④ Likelihood ratio exploits the following trick

gradient over policy function

why is this equivalent of the log?, just think of the chain rule

$$\begin{aligned}\nabla_\theta \pi_\theta(a|s) &= \pi_\theta(a|s) \frac{\nabla_\theta \pi_\theta(a|s)}{\pi_\theta(a|s)} && d[\log(f(x))] / dx = df(x) / f(x) \\ &= \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)\end{aligned}$$

this is like maximum likelihood estimation

- ⑤ The score function is $\nabla_\theta \log \pi_\theta(a|s)$ this is the gradient of our log likelihood function, where we use samples from MC
- ① Score indicates the steepness of the log-likelihood function
- ② [https://en.wikipedia.org/wiki/Score_\(statistics\)](https://en.wikipedia.org/wiki/Score_(statistics))

Policy Gradient for One-Step MDPs

- ① Consider a simple class of one-step MDPs

- ① Start in state $s \sim d(s)$ generate a starting s , based on the distribution of starting states
- ② Terminate after one time-step with reward $r = R(s, a)$ terminate immediately

- ② The objective is as follows

\theta are the parameters of the policy network

$$J(\theta) = \mathbb{E}_{\pi_\theta}[r]$$

this equation: we sum up the reward for every possible action in every possible state, weighted by how likely that action is going to happen $d(s)$ and how likely the agent is to pick the action $\pi_\theta(a|s)$

$$= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(a|s)r$$

REVIEW HOW THIS IS RELATED TO MLE

- ③ Use likelihood ratios to compute the policy gradient, then the gradient is as

taking the \grad, since its with respect to the policy, we can take it inside

$$\nabla_\theta J(\theta) = \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)r$$

then we use the MLE to get the log
this can be estimated through sampling

$$= \mathbb{E}_{\pi_\theta}[r \nabla_\theta \log \pi_\theta(a|s)]$$

expectation over cum. rewards

observations are state s

Discrete Policy Function: Softmax Policy

this slide is about how we use the features and weights to calculate a score we can softmax to get action probabilities

- ① Simple policy model: weight actions using a linear combination of features $\phi(s, a)^T \theta$
- ② Probability of action is proportional to the exponentiated weight

$$\pi_\theta(a|s) = \frac{\exp^{\phi(s,a)^T \theta}}{\sum_{a'} \exp^{\phi(s,a')^T \theta}}$$

how do we derive this?
its just basic calc and log rules

- ③ The score function is

something is a neural network?

$$\nabla_\theta \log \pi_\theta(a|s) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, .)]$$

features of action we actually took

average features of all actions, weighted by
how likely the policy thinks they are

Continuous Policy Function: Gaussian Policy

- ① In continuous action spaces, a Gaussian policy can be naturally defined
- ② Mean is a linear combination of state features $\mu(s) = \phi(s)^T \theta$
- ③ Variance may be fixed σ^2 or can also be parameterized
- ④ Policy is Gaussian value $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- ⑤ The score function is

$$\nabla_{\theta} \log \pi_{\theta}(a|s) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$

Policy Gradient for Multi-step MDPs

- ① Denote a sampled trajectory from π_θ from the environment as $\tau = (s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T, s_T) \sim (\pi_\theta, P(s_{t+1}|s_t, a_t))$
 - ① Sometimes we just simplify the notation as $\tau \sim \pi_\theta$
- ② Use $R(\tau) = \sum_{t=0}^T R(s_t, a_t)$ to be the sum of rewards for a trajectory τ
- ③ The policy objective is

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T R(s_t, a_t) \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

these are rollouts from the environment

where $P(\tau; \theta)$ denotes the probability over trajectories when executing the policy π_θ

- ④ Then our goal is to find the policy parameter θ

$$\theta^* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

Policy Gradient for Multi-step MDPs

- ① Our goal is to find the policy parameter θ

$$\theta^* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- ② Take the gradient with respect to θ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta)\end{aligned}$$

Policy Gradient for Multi-step MDPs

- ① Our goal is to find the policy parameter θ

$$\theta^* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- ② Take the gradient with respect to θ :

$$\nabla_{\theta} J(\theta) = \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta)$$

- ③ Approximate the gradient of objective with empirical estimate for m sample paths under policy π_{θ} :

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \nabla_{\theta} \log P(\tau_i; \theta)$$

$$\theta_{t+1} = \theta_t + \nabla_{\theta} J(\theta)$$

Decomposing a Trajectory

- ① Approximate with empirical estimate for m sample paths under policy π_θ :

$$\nabla_\theta J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \nabla_\theta \log P(\tau_i; \theta)$$

$\backslash\tau = \{s_0, a_0, r_1, s_1, \dots\}$

- ② Decompose $\nabla_\theta \log P(\tau; \theta)$

$\backslash\mu(s_0)$ is the initial distribution

$$\nabla_\theta \log P(\tau; \theta) = \nabla_\theta \log \left[\mu(s_0) \prod_{t=0}^{T-1} \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t) \right]$$

this is decided by the environment

might not be available

im guessing that $\log \backslash\mu(s_0)$ doesn't depend on $\backslash\theta$ so the gradient just takes it out

$$= \nabla_\theta \left[\log \mu(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t) \right]$$

this doesn't depend on the policy
so we can just take it out

$$= \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t)$$

Likelihood Ratio Policy Gradient

- ① Our goal is to find the policy parameter θ

$$\theta^* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- ② Approximate with empirical estimate for m sample paths under policy π_θ :

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \nabla_{\theta} \log P(\tau_i; \theta)$$

- ③ Thus we have $\nabla_{\theta} \log P(\tau_i; \theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)$$

- ① It does not need to know the transition dynamics!

Understand Score Function Gradient Estimator

- ① Consider the generic form of $E_{\tau \sim \pi_\theta}[R(\tau)]$ as

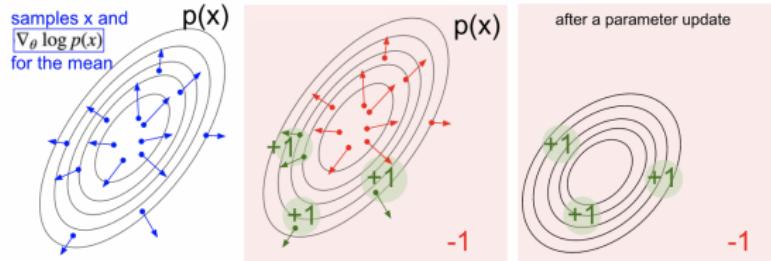
$$\nabla_\theta \mathbb{E}_{x \sim p(x; \theta)}[f(x)] = \mathbb{E}_x[f(x) \nabla_\theta \log p(x; \theta)]$$

x = trajectory

$p(x; \theta)$ = probability distribution of a trajectory happening

- ② The above gradient can be understood as:

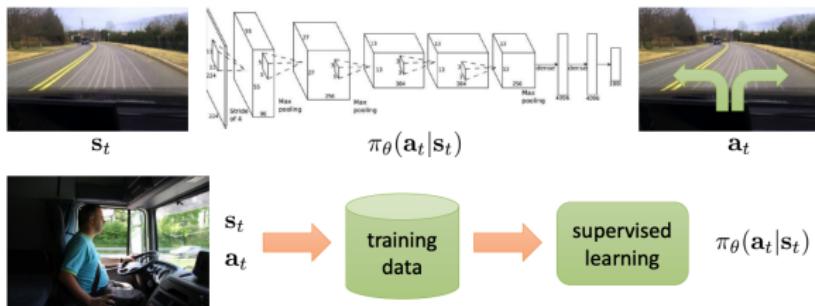
- ① Shift the distribution p through its parameter θ to let its future samples x achieve higher scores as judged by $f(x)$
- ② The direction of $f(x) \nabla_\theta \log p(x; \theta)$ pushes up the log probability of the sample, in proportion to how good it is



Comparison to Maximum Likelihood

- ① Policy gradient estimator: $\nabla_{\theta} J(\theta) \approx \frac{1}{M} \sum_{m=1}^M \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{t,m}, s_{t,m}) \right) \left(\sum_{t=1}^T r(s_{t,m}, a_{t,m}) \right)$

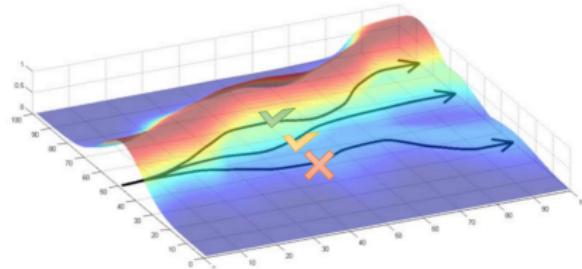
total reward for a trajectory
use a weighted sum of the gradient based on the reward
- ② Maximum likelihood estimator:
 $\nabla_{\theta} J_{ML}(\theta) \approx \frac{1}{M} \sum_{m=1}^M \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{t,m}, s_{t,m}) \right)$
- ③ good action is made more likely, bad action is made less likely



Comparison to Maximum Likelihood Estimation

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \nabla_{\theta} \log P(\tau_i; \theta)$$

- ① If going up the hill leads to higher reward, change the policy parameters to increase the likelihood of trajectories that move higher



Large Variance of Policy Gradient

- ① We have the following approximate update

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m R(\tau_i) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)$$

- ② Unbiased but having high variance as we use the MC method
- ③ Two fixes:
 - ① Use temporal causality
 - ② Include a baseline

Reduce Variance of Policy Gradient using Causality

in this, we go through an entire trajectory, then we use the score to reinforce every single action taken

this gives credit to actions that didn't actually cause the reward

- ① Previously $\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \right]$
- ② Causality: the policy cannot affect rewards in the past. This yields the following modified objective, where the sum of rewards here does not include the rewards achieved before the time step at which the policy is being queried: policy at time t' cannot affect reward at time t when $t < t'$

r1 is irrelevant to the actions and states $s_{\{1+}\), $a_{\{1+}\}$$

summation now starts at t and not 0

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{T-1} r_{t'} \right] \\ &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} G_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]\end{aligned}$$

ignore past rewards

we switch to Future rewards

Reduce Variance of Policy Gradient using Causality

- ① Therefore we have

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} G_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ② $G_t = \sum_{t'=t}^{T-1} r_{t'}$ is reward-to-go of a trajectory at step t , or a sample estimate of the Q function,
- ③ Then we can have the following estimated update

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{T-1} G_t^{(i)} \cdot \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)$$

REINFORCE: a Monte-Carlo policy gradient algorithm

- ① The algorithm simply samples multiple trajectories following the policy π_θ while updating θ using the estimated gradient

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

Input: a differentiable policy parameterization $\pi(a|s, \theta)$ define how we parameterize our policy

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$

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 $t = 0, \dots, T - 1$:

$G \leftarrow$ return from step t

go backwards because we are using the future rewards

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

- ② Classic paper: Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning: introduces REINFORCE algorithm

Reduce Variance Using a Baseline

① The original update

this is high variance because randomness in environment and policy

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} G_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ① $G_t = \sum_{t'=t}^{T-1} r_{t'}$ is the return for a trajectory which might have high variance
- ② We subtract a baseline b (that is a constant with respect to τ) from the policy gradient to reduce variance

mean of the gradient stays the same, but the variance is reduced
this is because we take the square distance from the mean and if the values are large,
square distance would be larger

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} (G_t - b) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ③ A good baseline can be the expected return

can approx this with monte carlo

$$b = \mathbb{E}[r_t + r_{t+1} + \dots + r_{T-1}]$$

Reduce Variance Using a Baseline

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} (\textcolor{red}{G_t - b}) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ① Interpretation: increase the logprob of action a_t proportional to how much returns G_t are better than the expected return
- ② We can **prove** that including baseline $b(s)$ can reduce variance, without changing the expectation:

$$\mathbb{E}_{\tau} \left[b \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] = 0, \quad (1)$$

$$\mathbb{E}_{\tau} \left[(G_t - b) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] = \mathbb{E}_{\tau} \left[G_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] \quad (2)$$

$$Var_{\tau} \left[(G_t - b) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] < Var_{\tau} \left[G_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] \quad (3)$$

- ③ Thus subtracting a baseline is unbiased in expectation but reduces variance of the sample gradient

Vanilla Policy Gradient Algorithm with Baseline

procedure POLICY GRADIENT(α)

 Initialize policy parameters θ and baseline values $b(s)$ for all s , e.g. to 0

for iteration = 1, 2, ... **do**

 Collect a set of m trajectories by executing the current policy π_θ

for each time step t of each trajectory $\tau^{(i)}$ **do**

 Compute the *return* $G_t^{(i)} = \sum_{t'=t}^{T-1} r_{t'}$

 Compute the *advantage estimate* $\hat{A}_t^{(i)} = G_t^{(i)} - b(s_t)$

 Re-fit the baseline to the empirical returns by updating w to minimize

this is the MSE loss function

$$\sum_{i=1}^m \sum_{t=0}^{T-1} \|b(s_t) - G_t^{(i)}\|^2$$

 Update policy parameters θ using the policy gradient estimate \hat{g}

$$\hat{g} = \sum_{i=1}^m \sum_{t=0}^{T-1} \hat{A}_t^{(i)} \nabla_\theta \log \pi_\theta(a_t^{(i)} | s_t^{(i)})$$

ghat is the estimated gradient

if we're using a softmax policy over some weights, gradient would be `gradient(theta = target action one hot encoded vector - probability vector`

 with an optimizer like SGD ($\theta \leftarrow \theta + \alpha \cdot \hat{g}$) or Adam

return θ and baseline values $b(s)$

- ① Classic paper: Sutton, McAllester, Singh, Mansour (1999). Policy gradient methods for reinforcement learning with function approximation: actor-critic algorithms with value function approximation

Practical Implementation of the Algorithm

- ① In practice we usually do not compute the gradients $\sum_t \hat{A}_t \nabla_\theta \log \pi_\theta(a_t|s_t)$ individually
- ② Instead, we accumulate data from current batch as

we can treat this as a "fake" objective function, so we get the policy gradient estimator when we differentiate

$$L(\theta) = \sum_t \hat{A}_t \log \pi_\theta(a_t|s_t; \theta)$$

- ③ Then the policy gradient estimator $\hat{g} = \nabla_\theta L(\theta)$
- ④ We also could have a joint loss with value function approximation as

$$L(\theta, \mathbf{w}) = \sum_t \left(\hat{A}_t \log \pi_\theta(a_t|s_t; \theta) - \min_{\mathbf{w}} \|b(s_t) - \hat{R}_t\|^2 \right)$$

- ⑤ Then solve this using auto diff

Rhat is the trajectory that we've calculated all (observed return)
w is used to approximate the second part of the loss

auto diff allows you to calculate loss without for loop

Reduce Variance Using a Critic

- ① The update is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} \textcolor{red}{G_t} \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ② In practice, G_t is a sample from Monte Carlo policy gradient, which is the unbiased but noisy estimate of $Q^{\pi_{\theta}}(s_t, a_t)$ q function is the expected return
- ③ Instead we can use a **critic** to estimate the action-value function,

$$Q_w(s, a) \approx Q^{\pi_{\theta}}(s, a)$$

- ④ Then the update becomes

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} \textcolor{red}{Q_w(s_t, a_t)} \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

Actor-Critic Policy Gradient

we used a learned prediction of the return
these use the weight w

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} Q_{\mathbf{w}}(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

This replaces G_t

- ① It becomes Actor-Critic Policy Gradient wants to maximize the value received from the critic
 - ① **Actor**: the policy function used to generate the action
 - ② **Critic**: the value function used to evaluate the reward of the actions
- ② Actor-critic algorithms maintain two sets of parameters
 - ① **Actor**: Updates policy parameters θ , in direction suggested by critic
 - ② **Critic**: Updates action-value function parameters \mathbf{w}

low variance and high bias if the critic is dumb

this way we don't need to finish an entire trajectory, we update on the go

Estimate the Action-Value Function

- ① Estimating the critic is solving a familiar problem: policy evaluation
 - ① How good is policy π_θ under current parameter θ
- ② Policy evaluation was explored in previous lectures, e.g.
 - ① Monte-Carlo policy evaluation
 - ② Temporal-Difference learning
 - ③ Least-squares policy evaluation

Action-Value Actor-Critic Algorithm

- ① Using a linear value function approximation: $Q_{\mathbf{w}}(s, a) = \psi(s, a)^T \mathbf{w}$
- ① Critic: update \mathbf{w} by a linear TD(0)
- ② Actor: update θ by policy gradient

features of (s, a)

linear approximator

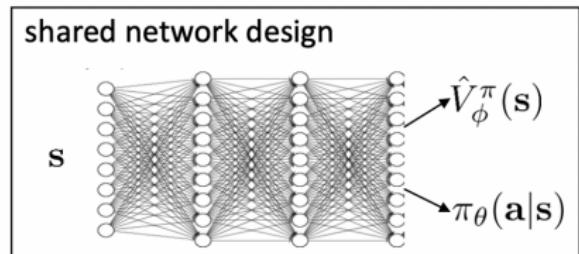
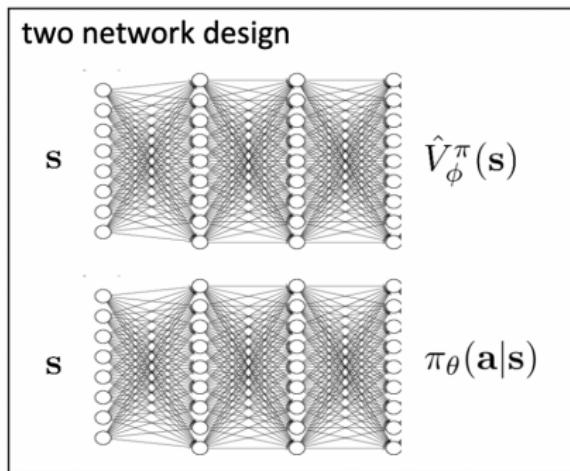
Algorithm 3 Simple QAC

```
1: for each step do
2:   generate sample  $s, a, r, s', a'$  following  $\pi_\theta$ 
3:    $\delta = r + \gamma Q_{\mathbf{w}}(s', a') - Q_{\mathbf{w}}(s, a)$  #TD error
4:    $\mathbf{w} \leftarrow \mathbf{w} - \beta \delta \psi(s, a)$  #Gradient descend
5:    $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(a|s) Q_{\mathbf{w}}(s, a)$  #Gradient ascend
6: end for
```

$\backslash\pi_\\theta$ is like a neural network where $\backslash\theta$ is the weights and biases, THIS MAKES SO MUCH MORE SENSE

Actor-Critic Function Approximators

- ① We can have two separate functions to approximate value function and policy function, or use a shared network design (feature extraction is shared but output two heads), as below:



Compatible Function Approximation

Theorem (Compatible Function Approximation Theorem)

If the following two conditions are satisfied:

1. *Value function approximator is compatible to the policy*

$$\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(a|s)$$

2. *Value function parameters w minimize the mean-squared error*

$$\epsilon(w) = E_{\pi_\theta}[(Q^{\pi_\theta}(s, a) - Q_w(s, a))^2]$$

then the policy gradient is exact

$$\nabla_\theta J(\theta) = E_{\pi_\theta}[Q_w(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

RS Sutton, et al. Policy gradient methods for reinforcement learning with function approximation. NIPS'01

Compatible Function Approximation

- ① Given a parameterized policy, compatible function approximation theorem can be used to derive an appropriate form of the value-function parameterization.
 - ① For example, consider a softmax policy as Gibbs distribution where ϕ_{sa} is the feature vector for state-action pair s,a

$$\pi_\theta(a|s) = \frac{e^{\theta^T \phi_{sa}}}{\sum_b e^{\theta^T \phi_{sb}}} \quad (4)$$

- ② Then meeting the compatibility condition requires

$$\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(a|s) = \phi_{sa} - \sum_b \pi_\theta(s, b) \phi_{sb}$$

- ③ so the compatible parameterization of $Q_w(s, a)$ is

$$Q_w(s, a) = w^T \left(\phi_{sa} - \sum_b \pi_\theta(s, b) \phi_{sb} \right) = w^T \nabla_\theta \log \pi_\theta(a|s)$$

Reducing the Variance of Actor-Critic by a Baseline

- ① Recall Q-function / state-action-value function:

(discounted future rewards)

$$Q^{\pi, \gamma}(s, a) = \mathbb{E}_{\pi}[r_1 + \gamma r_2 + \dots | s_1 = s, a_1 = a]$$

- ② State value function can serve as a great baseline

averaged over all the actions at a state

(weighted by the probability that we choose an action)

$$\begin{aligned} V^{\pi, \gamma}(s) &= \mathbb{E}_{\pi}[r_1 + \gamma r_2 + \dots | s_1 = s] \\ &= \mathbb{E}_{a \sim \pi}[Q^{\pi, \gamma}(s, a)] \end{aligned}$$

- ③ Advantage function: combining Q with baseline V

$$A^{\pi, \gamma}(s, a) = Q^{\pi, \gamma}(s, a) - V^{\pi, \gamma}(s)$$

- ④ Then the policy gradient becomes:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi, \gamma}(s, a)]$$

Advantage Actor-Critic

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi_{\theta}}(s, a)]$$

- ① Critic should estimate the advantage function
- ② Use two function approximators and two sets of parameter vectors,

$$V_v(s) \approx V^{\pi}(s)$$
$$Q_w(s, a) \approx Q^{\pi}(s, a)$$

- ③ Update both parameters by TD learning or MC

Advantage Actor-Critic

- ① For the true value function $V^{\pi_\theta}(s)$, the TD error δ^{π_θ}
discounted value of the next state

$$\delta^{\pi_\theta} = r(s, a) + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$$

baseline

is an unbiased estimate of the advantage function

$$\begin{aligned}\mathbb{E}_{\pi_\theta}[\delta^{\pi_\theta}|s, a] &= \mathbb{E}_{\pi_\theta}[r + \gamma V^{\pi_\theta}(s')|s, a] - V^{\pi_\theta}(s) \\ &= Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) \\ &= A^{\pi_\theta}(s, a)\end{aligned}$$

- ② So we can use the TD error to compute the policy gradient

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s) \delta^{\pi_\theta}]$$

- ③ It requires only **one set of critic parameter κ** that is estimated by TD
since we don't have a q value function, we just need to have one set of weights for the value function

$$\delta_V = r + \gamma V_\kappa(s') - V_\kappa(s)$$

N-step estimators

- ① We used the Monte-Carlo estimates of the reward
- ② We can also use TD methods for the policy gradient update, or any intermediate blend between TD and MC methods:
- ③ Consider the following n -step returns for $n = 1, 2, \infty$

$$n = 1(TD) \quad G_t^{(1)} = r_{t+1} + \gamma v(s_{t+1})$$

$$n = 2 \quad G_t^{(2)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2})$$

$$n = \infty(MC) \quad G_t^{(\infty)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-t-1} r_T$$

- ④ Then the advantage estimators become

$$\hat{A}_t^{(1)} = r_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\hat{A}_t^{(2)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2}) - v(s_t)$$

$$\hat{A}_t^{(\infty)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-t-1} r_T - v(s_t)$$

- ⑤ $\hat{A}^{(1)}$ has low variance and high bias. $\hat{A}^{(\infty)}$ has high variance but low bias

Critic at Different Time-Scales

- ① Critic can estimate linear value function $V_\kappa(s) = \psi(s)^T \kappa$ from many targets as follows
- ② For MC, the update is

$$\Delta \kappa = \alpha(\textcolor{red}{G_t} - V_\kappa(s))\psi(s)$$

- ③ For TD(0), the update is

$$\Delta \kappa = \alpha(\textcolor{red}{r + \gamma V_\kappa(s')} - V_\kappa(s))\psi(s)$$

- ④ For k-step return, the update is

$$\Delta \kappa = \alpha \left(\sum_{i=0}^k \textcolor{red}{\gamma^i r_{t+i}} + \textcolor{red}{\gamma^k V_\kappa(s_{t+k})} - V_\kappa(s_t) \right) \psi(s)$$

Actors at Different Time-Scales

- ① The policy gradient can also be estimated at many time-scales

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi_{\theta}}(s, a)]$$

- ② Monte-Carlo Actor-Critic policy gradient uses error from complete return

$$\Delta\theta = \alpha(G_t - V_{\kappa}(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- ③ TD Difference Actor-Critic policy gradient uses TD error

$$\Delta\theta = \alpha(r + \gamma V_{\kappa}(s_{t+1}) - V_{\kappa}(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- ④ k-step return Actor-Critic policy gradient uses the k-step return error

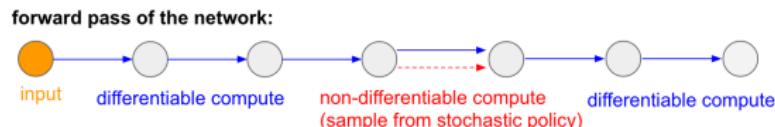
$$\Delta\theta = \alpha \left(\sum_{i=0}^k \gamma^i r_{t+i} + \gamma^k V_{\kappa}(s_{t+k}) - V_{\kappa}(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

Property of PG: Policy Gradient is On-Policy RL

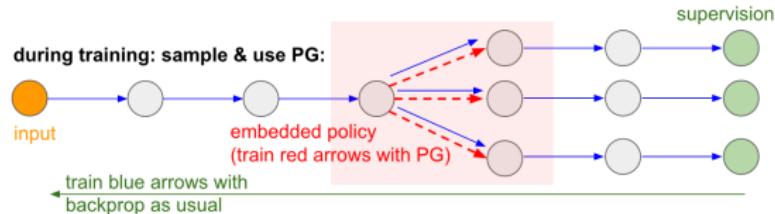
- ① $\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P_{\theta}(\tau) r(\tau)]$
- ② In REINFORCE algorithm: there is the on-policy sampling (sample $\{\tau\}$ from π_{θ})
- ③ On-policy learning is stable but might be inefficient

Property of PG: Overcoming Non-differentiable Computation

- ① Another interesting strength of Policy Gradient is that it allows us to overcome the non-differentiable computation



- ② During training we will produce several samples (indicated by the branches below), and then we'll encourage samples that eventually led to good outcomes (in this case for example measured by the loss at the end)



Summary of Policy Gradient

- ① Softmax policy: weight actions using linear combination of features
 $\phi(s, a)^T \theta$

$$\pi_\theta(a|s) = \frac{\exp^{\phi(s,a)^T \theta}}{\sum_{a'} \exp^{\phi(s,a')^T \theta}}$$

- ② Gaussian policy:
 - ① Mean is a linear combination of state features $\mu(s) = \phi(s)^T \theta$
 - ② Variance may be fixed σ^2 or can also be parameterized
 - ③ Policy is Gaussian, the continuous $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- ③ Neural network policy: $\pi_\theta(a|s) = F(s)$ where $F(\cdot)$ is a neural network

Summary of Policy Gradient

- ① In policy optimization, for the policy function π_θ the objective is to maximize

this is an on-policy loss function

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

- ① $R(\tau)$ could be any reasonable reward function on τ
- ② So the gradient is

$$\begin{aligned}\nabla_\theta J(\theta) &= \nabla_\theta \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \mathbb{E}_\tau [R(\tau) \nabla_\theta \log P(\tau; \theta)]\end{aligned}$$

- ① **log likelihood trick:** $\nabla_\theta f_\theta(x) = f_\theta(x) \nabla_\theta \log f_\theta(x)$
- ② Trajectory distribution: $P(\tau; \theta) = \mu(s_0) \prod_{t=0}^{T-1} \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$
- ③ **cancelling trick:** $\nabla_\theta \log P(\tau; \theta) = \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t)$

Summary of Policy Gradient

- ① After decomposing τ we derived that

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \right]$$

- ② After applying the causality, we have

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} G_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ① $G_t = \sum_{t'=t}^{T-1} r_{t'}$ is the return for a trajectory at step t

Summary of Policy Gradient

- ① The policy gradient has many forms

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau}[\textcolor{red}{G_t} \nabla_{\theta} \log \pi_{\theta}(a|s)] \text{ --- REINFORCE}$$

$$= \mathbb{E}_{\tau}[\textcolor{red}{Q^{\pi}(s, a)} \nabla_{\theta} \log \pi_{\theta}(a|s)] \text{ --- Q Actor-Critic}$$

$$= \mathbb{E}_{\tau}[\textcolor{red}{A^{\pi}(s, a)} \nabla_{\theta} \log \pi_{\theta}(a|s)] \text{ --- Advantage Actor-Critic}$$

$$= \mathbb{E}_{\tau}[\delta \nabla_{\theta} \log \pi_{\theta}(a|s)] \text{ --- TD Actor-Critic}$$

- ② We can use policy evaluation (e.g. MC or TD learning) to estimate $Q^{\pi}(s, a)$, $A^{\pi}(s, a)$, or $V^{\pi}(s, a)$

Extension of Policy Gradient

- ① State-of-the-art RL methods are almost all policy-based
 - ① **A2C and A3C:** Asynchronous Methods for Deep Reinforcement Learning, ICML'16. Representative high-performance actor-critic algorithm: <https://openai.com/blog/baselines-acktr-a2c/>
 - ② **TRPO:** Schulman, L., Moritz, Jordan, Abbeel (2015). Trust region policy optimization: deep RL with natural policy gradient and adaptive step size
 - ③ **PPO:** Schulman, Wolski, Dhariwal, Radford, Klimov (2017). Proximal policy optimization algorithms: deep RL with importance sampled policy gradient

Resource on policy gradient

- ① A very nice summary of policy gradient algorithms:

<https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html>

- ② Educational blog by Karpathy:

<http://karpathy.github.io/2016/05/31/r1/>

Code example

- ① REINFORCE code on CartPole:

<https://github.com/ucla-rlcourse/RExample/blob/master/policygradient/reinforce.py>

- ② Policy Gradient on Pong : <https://github.com/ucla-rlcourse/RExample/blob/master/policygradient/pg-pong-pytorch.py>

- ③ Policy Gradient with Baseline on Pong:

<https://github.com/ucla-rlcourse/RExample/blob/master/policygradient/pgb-pong-pytorch.py>

- ④ Actor Critic on Pong: <https://github.com/ucla-rlcourse/RExample/blob/master/policygradient/ac-pong-pytorch.py>