

Policy-Gradient Methods

Ayça Özçelikkale

Dept. of Electrical Engineering, Uppsala University

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Overview of the Course

- ▶ Planning by Dynamic Programming
- ▶ Model-free prediction and control with tabular methods
- ▶ Function Approximation: We can use function approximation
 - ▶ to approximate value functions (next lecture in Actor-Critic)
 - ▶ to parametrize policies (this lecture and the next lecture)
 - ▶ to model the environment[†]
- ▶ Model-based RL
 - ▶ Tabular Model-based RL Methods: Dyna-Q and Dyna-Q+ (last lecture)
 - ▶ General Model-based RL Methods:[†]
- ▶ Policy-Gradient Methods (this lecture and the next lecture)
 - ▶ REINFORCE
 - ▶ REINFORCE with baseline
 - ▶ One-step Actor Critic

[†] 5hp course: You're not responsible to know the details of algorithms/methods but you should be aware of the general ideas and advantages/disadvantages

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

MAIN IDEA: The policy determines which actions to take. Let's learn a policy directly instead of relying on value functions!

What are Policy-Gradient Methods?

Short Review of Previous Lectures: Up to now, we have used the value or the action-value function to find a good policy. Then, a policy is generated directly from the value functions.

► **Example:** Q-learning

- The learned action-value function Q approximates the optimal action-value function $q_*(s, a)$
- The policy is implicit, for instance we use ϵ -greedy.

Now: We will parametrize the policy with $\theta \in \mathbb{R}^{d'}$:

$$\pi_{\theta}(a|s) = P[A_t = a | S_t = s, \theta_t = \theta]$$

$\pi_{\theta}(a|s)$: Probability action a is taken at time t given that the environment is in state s and the policy is parametrized by the parameter θ .

We will now directly search for a good policy π by searching over the parameters θ .

Example - Policy Parametrization

Let the action space be discrete. We can use **soft-max** for policy parametrization:

- ▶ $h_{\theta}(s, a)$ denotes preferences: $h_{\theta}(s, a)$ is high when action a is preferred over other actions)
- ▶ Action probabilities are proportional to exponentiated preferences:

$$\pi_{\theta}(a|s) \propto e^{h_{\theta}(s, a)}$$

In particular, set

$$\pi_{\theta}(a|s) \triangleq \frac{e^{h_{\theta}(s, a)}}{\sum_b e^{h_{\theta}(s, b)}}$$

The above scheme is called soft-max in action preferences.

How to parametrize $h_{\theta}(s, a)$?

- ▶ Example: Linear in the features $\mathbf{x}(s, a)$

$$h_{\theta}(s, a) = \theta^T \mathbf{x}(s, a)$$

- ▶ Example: $\text{ANN}_{\theta}(s, a)$

Why are these methods called Policy Gradient Methods?

Let $J(\theta)$ be our scalar performance measure. Finding the best policy is now equivalent to maximizing $J(\theta)$ over θ :

$$\max_{\theta \in \mathbb{R}^{d'}} J(\theta)$$

Policy gradient methods seek for θ using (approximate) gradient ascent. Hence, the updates for policy gradient methods approximate the **gradient ascent** updates

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$$

where $\alpha > 0$ is the step size parameter.

Classification of RL Approaches

Value-Based

- ▶ Learn Value Function
- ▶ Policy is Implicit

Actor-Critic

- ▶ Learn Value Function (Critic)
- ▶ Learn Policy (Actor)

Policy-Based

- ▶ No Value Function
- ▶ Learn Policy

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Motivation

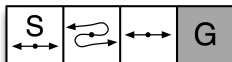
Policy-gradient methods are promising from the following perspectives:

- ▶ Deal with high-dimensional or continuous action spaces directly
- ▶ Directly learn stochastic policies
 - ▶ In some environments, we absolutely need stochastic policies
 - ▶ Example: Poker
- ▶ Better convergence properties
 - ▶ Action probabilities change smoothly as a function of the policy parameters → It is possible to directly approximate gradient ascent on the performance
- ▶ Incorporate a prior information about the problem

Note that policy-gradient methods also have some disadvantages:

- ▶ Typically converge to a local rather than global optimum
- ▶ Evaluating a policy is typically inefficient and high variance

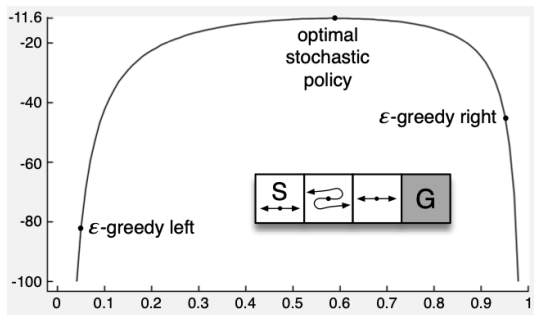
Example: Short corridor with switched actions



This is an example where the optimal policy is stochastic.

- ▶ A corridor of three rooms and one terminating state at the right.
 - ▶ We start at the left-most room S .
- ▶ Reward: -1 per step
- ▶ Actions: Left, Right
 - ▶ Middle room reverses the actions!
- ▶ Features: $x(s, right) = [1, 0]^T$ and $x(s, left) = [0, 1]^T$, $\forall s$
 - ▶ All states appear identical under function approximation!
- ▶ An action-value method with ϵ -greedy can only have two forms:
 - ▶ Right with probability $1 - \epsilon/2$, $\forall s$
 - ▶ Left with probability $1 - \epsilon/2$, $\forall s$

Optimal policy is stochastic



Performance $J(\theta) = v_{\pi_\theta}(S)$ versus the probability of choosing action "right"

For ϵ -greedy left/right, $\epsilon = 0.1$

Exercise: Action Values vs Action Preferences

Consider the policy parametrization with soft-max in action preferences:

$$\pi_{\theta}(a|s) \triangleq \frac{e^{h_{\theta}(s,a)}}{\sum_b e^{h_{\theta}(s,b)}}$$

An idea is to use to first estimate the action values $q(s, a)$ with a value-based method; and then use $\hat{q}(s, a)$ instead of $h_{\theta}(s, a)$ in the above equation.

Do you think that this idea will work?

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Reminder: Gradient Ascent -A general idea for minimizing functions

Notation- Gradient: Let $J(\theta) : \mathbb{R}^{d'} \rightarrow \mathbb{R}$ be a scalar valued function with a vector input. Gradient of $J(\theta)$ is given by

$$\nabla J \triangleq \frac{\partial J}{\partial \theta} \triangleq \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_{d'}} \end{bmatrix} \in \mathbb{R}^{d' \times 1},$$

The gradient evaluated at a particular value $\bar{\theta}$ is denoted by the following:

$$\nabla J(\bar{\theta}) = \left. \frac{\partial J}{\partial \theta} \right|_{\theta=\bar{\theta}}$$

Gradient Ascent Idea: If you want to minimize a function $J(\theta)$ iteratively, start from an initial point θ_i (here i is the iteration index!), and move in the direction of maximum “ascent”:

$$\begin{aligned} \theta_{i+1} &= \theta_i + \Delta \theta_i \\ &= \theta_i + \alpha^\theta \nabla J(\theta_i) \end{aligned}$$

where $\alpha^\theta > 0$ is the step size.

Objective Function $J(\theta)$

- ▶ We focus on the episodic case
- ▶ We assume that every episode starts at some particular non-random state s_0
- ▶ The performance measure is defined as the value function evaluated at the starting state s_0 if we follow the policy π_θ

$$J(\theta) \triangleq v_{\pi_\theta}(s_0)$$

- ▶ For the sake of clarity, analytical development assumes no discounting, i.e. $\gamma = 1$

How to find the policy gradient?

Challenge: Policy parameter θ affects both of the following

- ▶ Action selections
 - ▶ Given a state, the effect of θ on the actions and the reward can be calculated easily
- ▶ State distribution
 - ▶ The effect of θ on the state distribution is a function of the unknown environment

Policy Gradient Theorem gives a neat answer to this challenge:

$$\nabla J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_a q_{\pi_{\theta}}(S_t, a) \nabla \pi_{\theta}(a|S_t) \right]$$

- ▶ The expectation is over the state distribution under π_{θ} (We know how to approximate this type of expressions!)
- ▶ RHS does not involve the derivative of the state distribution!

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

- REINFORCE update

- Algorithm

- Example

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE update

Algorithm

Example

REINFORCE with Baseline

Actor-Critic

One-step Actor-Critic

Interpretation in Terms of Advantage Function

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

We look at $\nabla J(\boldsymbol{\theta})$ more carefully:

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[\sum_a q_{\pi_{\boldsymbol{\theta}}}(S_t, a) \nabla \pi_{\boldsymbol{\theta}}(a|S_t) \right] \\&= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[\sum_a \pi_{\boldsymbol{\theta}}(a|S_t) q_{\pi_{\boldsymbol{\theta}}}(S_t, a) \frac{\nabla \pi_{\boldsymbol{\theta}}(a|S_t)}{\pi_{\boldsymbol{\theta}}(a|S_t)} \right] \\&= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[q_{\pi_{\boldsymbol{\theta}}}(S_t, A_t) \frac{\nabla \pi_{\boldsymbol{\theta}}(A_t|S_t)}{\pi_{\boldsymbol{\theta}}(A_t|S_t)} \right] \\&= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[G_t \frac{\nabla \pi_{\boldsymbol{\theta}}(A_t|S_t)}{\pi_{\boldsymbol{\theta}}(A_t|S_t)} \right] \\&= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} [G_t \nabla \ln \pi_{\boldsymbol{\theta}}(A_t|S_t)]\end{aligned}$$

where in the last two lines we have used $\mathbb{E}_{\pi}[G_t|S_t, A_t] = q_{\pi}(S_t, A_t)$ and $\nabla \ln f(\boldsymbol{\theta}) = \frac{\nabla f(\boldsymbol{\theta})}{f(\boldsymbol{\theta})}$

We use $\nabla J(\theta)$ for stochastic gradient ascent updates

Gradient Ascent: To maximize $J(\theta)$, start from an estimate θ_t and update using

$$\theta_{t+1} = \theta_t + \alpha^\theta \nabla J(\theta_t)$$

From the previous slide: $\nabla J(\theta) = \mathbb{E}_{\pi_\theta} [G_t \nabla \ln \pi_\theta(A_t|S_t)]$

Hence, we have

$$\theta_{t+1} = \theta_t + \alpha^\theta \mathbb{E}_{\pi_\theta} [G_t \nabla \ln \pi_{\theta_t}(A_t|S_t)]$$

where $\nabla \ln \pi_{\theta_t}(\cdot)$ is the gradient of $\ln \pi_\theta(\cdot)$ with respect to θ evaluated at the particular value θ_t .

Stochastic Gradient Ascent: Use a sample instead of an expectation

$$\theta_{t+1} = \theta_t + \alpha^\theta G_t \nabla \ln \pi_{\theta_t}(A_t|S_t)$$

This is the **REINFORCE** update.

We look at REINFORCE updates more closely

REINFORCE update from the previous slide:

$$\begin{aligned}\theta_{t+1} &= \theta_t + \Delta\theta_t \\ &= \theta_t + \alpha^\theta G_t \nabla \ln \pi_{\theta_t}(A_t|S_t) \\ &= \theta_t + \alpha^\theta G_t \frac{\nabla \pi_{\theta_t}(A_t|S_t)}{\pi_{\theta_t}(A_t|S_t)}\end{aligned}$$

Note that $\nabla \pi_{\theta_t}(A_t|S_t)$ is the direction in θ space that most increases the probability of taking action A_t when at state S_t

Each update $\Delta\theta_t$ is in the direction of $\nabla \pi_{\theta_t}(A_t|S_t)$ with a scaling so that

- ▶ Proportional to the return G_t
- ▶ Inversely proportional to the action probability

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE update

Algorithm

Example

REINFORCE with Baseline

Actor-Critic

One-step Actor-Critic

Interpretation in Terms of Advantage Function

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Reinforce: Monte-Carlo Policy Gradient Control for finding π_*

Input: Differentiable policy parametrization $\pi_{\theta}(a|s)$.

Initialization: Set step size $\alpha^{\theta} > 0$. Initialize the policy parameters $\theta \in \mathbb{R}^{d'}$.

repeat (for each episode)

Generate an episode using π_{θ} : $S_0, A_0, R_1, S_1, A_1, \dots, R_T$.

repeat : (for each step of the episode $t = 0, 1, \dots, T-1$)

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\theta \leftarrow \theta + \alpha^{\theta} \gamma^t G_t \nabla \ln \pi_{\theta}(A_t | S_t)$$

Desired Output $\pi_{\theta} \approx \pi_*$

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE update

Algorithm

Example

REINFORCE with Baseline

Actor-Critic

One-step Actor-Critic

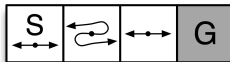
Interpretation in Terms of Advantage Function

Policy Parametrization for Continuous Actions

Discussions

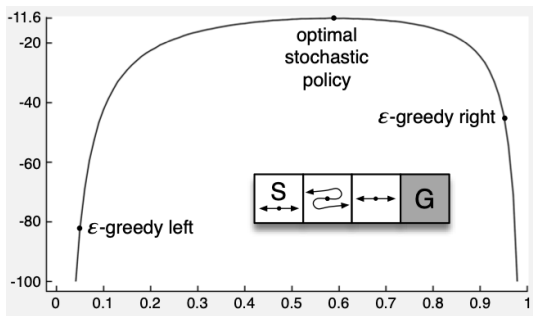
Concluding Remarks

Recall the environment “Short corridor with switched actions”



- ▶ A corridor of three rooms and one terminating state at the right.
 - ▶ We start at the left-most room S .
- ▶ Reward: -1 per step
- ▶ Actions: Left, Right
 - ▶ Middle room reverses the actions!
- ▶ Features: $x(s, right) = [1, 0]^T$ and $x(s, left) = [0, 1]^T, \forall s$
 - ▶ All states appear identical under function approximation!
- ▶ An action-value method with ϵ -greedy can only have two forms:
 - ▶ Right with probability $1 - \epsilon/2, \forall s$
 - ▶ Left with probability $1 - \epsilon/2, \forall s$

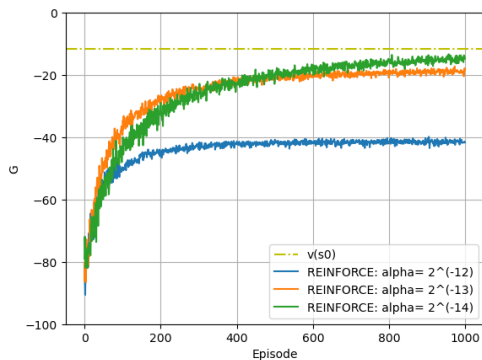
Reminder: Optimal policy is stochastic



Performance $J(\theta) = v_{\pi_\theta}(S)$ versus the probability of choosing action "right"

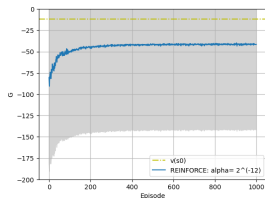
For ϵ -greedy left/right, $\epsilon = 0.1$

Convergence of REINFORCE

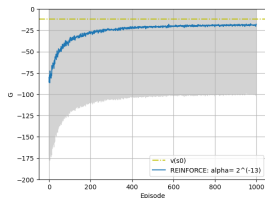


Total reward on episode (Average over 300 runs)

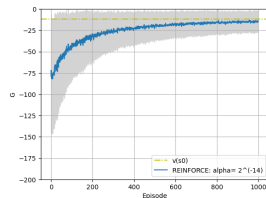
REINFORCE is of high variance



$$\alpha^{\theta} = 2^{-12}$$



$$\alpha^{\theta} = 2^{-13}$$

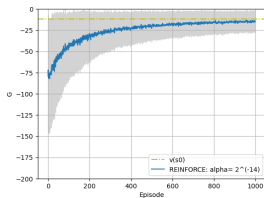


$$\alpha^{\theta} = 2^{-14}$$

- ▶ Blue curves: mean of the total reward over runs ,i.e. same with the corresponding curve in the former figure
- ▶ We color the area between “mean-standard deviation” and “mean+standard deviation” with light gray
- ▶ To make the plots easier to read, we have done the following: standard deviation is smoothed along episode axis and capped with 100

REINFORCE with baseline has lower variance

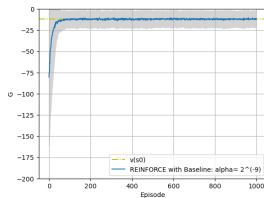
REINFORCE



$$\alpha^{\theta} = 2^{-14}$$

V.S.

REINFORCE with baseline



$$\alpha^{\theta} = 2^{-9}$$

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Reminders: Classification of RL Approaches

Value-Based

- ▶ Learn Value Function
- ▶ Policy is Implicit

Actor-Critic

- ▶ Learn Value Function (Critic)
- ▶ Learn Policy (Actor)

Policy-Based

- ▶ No Value Function
- ▶ Learn Policy

Reminders:

Policy $\pi_{\theta}(a|s)$ is parametrized by the policy parameter vector $\theta \in \mathbb{R}^{d'}$

Performance measure is value function evaluated at the starting state s_0 if we follow the policy π_{θ} :

$$J(\theta) \triangleq v_{\pi_{\theta}}(s_0)$$

Aim: Find the best policy by maximizing $J(\theta)$ over θ : $\max_{\theta \in \mathbb{R}^{d'}} J(\theta)$

Main Idea for Policy Gradient Methods:

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$$

Policy Gradient Theorem:

$$\nabla J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_a q_{\pi_{\theta}}(S_t, a) \nabla \pi_{\theta}(a|S_t) \right] = \mathbb{E}_{\pi_{\theta}} [G_t \nabla \ln \pi_{\theta}(A_t|S_t)]$$

REINFORCE Update:

$$\theta \leftarrow \theta + \alpha G_t \nabla \ln \pi_{\theta}(A_t|S_t)$$

Preliminaries: Baseline[†]

We can add a baseline $b(s)$ to obtain a generalization of the policy-gradient theorem

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\sum_a q_{\pi_{\theta}}(S_t, a) \nabla \pi_{\theta}(a|S_t) \right] \\ &= \mathbb{E}_{\pi_{\theta}} \left[\sum_a (q_{\pi_{\theta}}(S_t, a) - b(S_t)) \nabla \pi_{\theta}(a|S_t) \right]\end{aligned}$$

since we have

$$\begin{aligned}\sum_a b(s) \nabla \pi_{\theta}(a|s) &= b(s) \sum_a \nabla \pi_{\theta}(a|s) = b(s) \nabla \sum_a \pi_{\theta}(a|s) \\ &= b(s) \nabla 1 \\ &= 0\end{aligned}$$

- ▶ This equation holds as long as the baseline does not depend on the actions.
- ▶ Note that the baseline can depend on the state s

[†] : Definition of baseline from Merriam-Webster dictionary: a line serving as a basis, especially one of known measure or position used (as in surveying or navigation) to calculate or locate something

REINFORCE with baseline

Update for REINFORCE with baseline

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} (G_t - b(S_t)) \nabla \ln \pi_{\boldsymbol{\theta}}(A_t | S_t)$$

- ▶ Baseline does not change the expected value of the update
- ▶ For $b(s) = 0$, we obtain the plain REINFORCE algorithm.
- ▶ The baseline $b(s)$ should vary with the state.
 - ▶ Example: If all actions from a given state has high values, it is better to have a high baseline so that we can differentiate between the actions.
- ▶ **Idea:** Use $b(s) = \hat{v}(S_t, \mathbf{w})$ where $\hat{v}(\cdot)$ is an estimate parametrized by $\mathbf{w} \in \mathbb{R}^d$.

Reinforce with Baseline for finding π_*

Input: Differentiable policy parametrization $\pi_{\theta}(a|s)$; differentiable value function parametrization $\hat{v}(s, \mathbf{w})$

Initialization: Set step sizes $\alpha^{\theta} > 0$, $\alpha^{\mathbf{w}} > 0$. Initialize the policy parameters $\theta \in \mathbb{R}^{d'}$ and value function approximation parameters $\mathbf{w} \in \mathbb{R}^d$.

repeat (for each episode)

Generate an episode using π_{θ} : $S_0, A_0, R_1, S_1, A_1, \dots, R_T$.

repeat : (for each step of the episode $t = 0, 1, \dots, T-1$)

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\delta = G - \hat{v}(S_t, \mathbf{w})$$

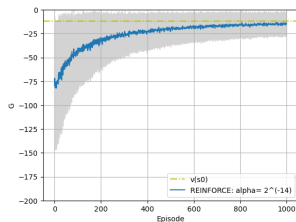
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$$

$$\theta \leftarrow \theta + \alpha^{\theta} \gamma^t \delta \nabla \ln \pi_{\theta}(A_t | S_t)$$

Desired Output $\pi_{\theta} \approx \pi_*$

Example Revisited: Short Corridor with switched actions

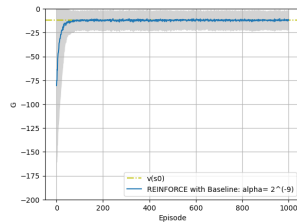
REINFORCE



$$\alpha^\theta = 2^{-14}$$

V.S.

REINFORCE with baseline



$$\alpha^\theta = 2^{-9}, \alpha^w = 2^{-6}$$

- ▶ REINFORCE with baseline has lower variance.
- ▶ REINFORCE with baseline learns faster.
- ▶ How to set step sizes α^w and α^θ ?
 - ▶ α^w : We have a rule of thumb

$$\alpha^w = \frac{0.1}{\mathbb{E}[\|\nabla \hat{v}(S_t, \mathbf{w})\|^2]} = \frac{0.1}{\mathbb{E}[\|\mathbf{x}(S_t, \mathbf{w})\|^2]}$$

- ▶ α^θ : More difficult to choose.

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

- One-step Actor-Critic

- Interpretation in Terms of Advantage Function

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

- REINFORCE update

- Algorithm

- Example

REINFORCE with Baseline

Actor-Critic

- One-step Actor-Critic**

- Interpretation in Terms of Advantage Function

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Actor-Critic Methods

Actor-critic methods is a large family of methods which uses the following **main idea**:
Use the value function not only as a baseline but also as a “critic”

- critic uses the value function for bootstrapping, updating the value estimate for a state from the estimated values of subsequent states

Note: Although the conceptual distinction is important, algorithms with a state-dependent baseline is sometimes classified also as actor-critic algorithms.

Example: One-step Actor-Critic

REINFORCE with baseline

$$\theta \leftarrow \theta + \alpha^{\theta} (G_t - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_{\theta}(A_t | S_t)$$

One-step Actor-Critic

$$\theta \leftarrow \theta + \alpha^{\theta} (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_{\theta}(A_t | S_t)$$

One-step Actor-Critic for finding π_*

Input: Differentiable policy parametrization $\pi_{\theta}(a|s)$; differentiable value function parametrization $\hat{v}(s, \mathbf{w})$

Initialization: Set step sizes $\alpha^{\theta} > 0$, $\alpha^{\mathbf{w}} > 0$. Initialize the policy parameters $\theta \in \mathbb{R}^{d'}$ and value function approximation parameters $\mathbf{w} \in \mathbb{R}^d$.
repeat (for each episode)

Initialize first state S of the episode.

$l \leftarrow 1$

repeat : (for each time step, while S is not terminal)

Choose action A using $\pi_{\theta}(\cdot|S)$

Take action A , observe S' , R

$\delta = R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$ (If S' terminal, set $\hat{v}(S', \mathbf{w})$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$

$\theta \leftarrow \theta + \alpha^{\theta} \delta \nabla \ln \pi_{\theta}(A|S)$

$l \leftarrow \gamma l$

$S \leftarrow S'$

Desired Output $\pi_{\theta} \approx \pi_*$

REINFORCE

$$\theta \leftarrow \theta + \alpha^\theta G_t \nabla \ln \pi_\theta(A_t | S_t)$$

- high variance, updates unbiased, inconvenient for online/continuing problems

REINFORCE with a state-dependent baseline $\hat{v}(S_t, \mathbf{w})$ found using G_t

$$\theta \leftarrow \theta + \alpha^\theta (G_t - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_\theta(A_t | S_t)$$

- lower variance, updates unbiased, inconvenient for online/continuing problems

One-step Actor-Critic

$$\theta \leftarrow \theta + \alpha^\theta (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_\theta(A_t | S_t)$$

- lower variance, updates biased, convenient for online/continuing problems

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

- REINFORCE update

- Algorithm

- Example

REINFORCE with Baseline

Actor-Critic

- One-step Actor-Critic

- Interpretation in Terms of Advantage Function**

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Interpreting One-step Actor-Critic in terms of Advantage Function

Advantage Function:

$$A_{\pi_{\theta}}(s, a) = q_{\pi_{\theta}}(s, a) - v_{\pi_{\theta}}(s)$$

We can write policy gradient in terms of the advantage function

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\pi_{\theta}} [q_{\pi_{\theta}}(S_t, A_t) \nabla \ln \pi_{\theta}(A_t | S_t)] \\ &= \mathbb{E}_{\pi_{\theta}} [A_{\pi_{\theta}}(S_t, A_t) \nabla \ln \pi_{\theta}(A_t | S_t)]\end{aligned}$$

REINFORCE: Approximate policy gradients by estimating $q_{\pi_{\theta}}(S_t, A_t)$

Actor-Critic: Approximate policy gradients by estimating $A_{\pi_{\theta}}(S_t, A_t)$

- ▶ Example: Use function approximation for both $q_{\pi_{\theta}}(s, a)$ and $v_{\pi_{\theta}}(s)$
- ▶ Example (One-Step Actor-Critic): Use function approximation for $v_{\pi_{\theta}}(s)$ and use TD idea to find estimates for $q_{\pi_{\theta}}(s, a)$

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Preliminaries: Normal Distribution

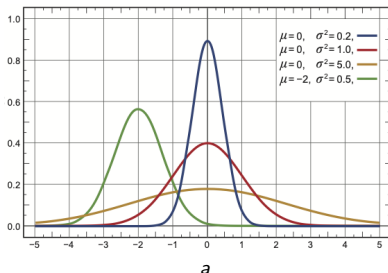
Probability density function (pdf) of a scalar random variable x is defined by the following relationship:

$$\text{Prob}[z_L \leq a \leq z_H] = \int_{z_L}^{z_H} p_a(\bar{a}) d\bar{a}$$

The pdf for normal (i.e. Gaussian) distribution is given by¹

$$p_a(a) = \mathcal{N}(a; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(a-\mu)^2}{2\sigma^2}\right)$$

where μ is the mean and σ is the standard deviation.



¹ π does not denote the policy in this formula; it is just the number $\pi \approx 3.1415$

Policy Parametrization for Continuous Actions

Popular Idea: Define the policy π as the normal probability density over a real-valued scalar action

$$\pi(a|s, \theta) = \mathcal{N}(a; \mu(s, \theta), \sigma(s, \theta))$$

where $\mu(s, \theta)$ and $\sigma(s, \theta)$ are function approximators parametrized by θ .

Example:

$$\theta = \begin{bmatrix} \theta_\mu \\ \theta_\sigma \end{bmatrix}$$

$$\mu(s, \theta) = \theta_\mu^T \mathbf{x}_\mu(s)$$

$$\sigma(s, \theta) = \exp(\theta_\sigma^T \mathbf{x}_\sigma(s))$$

where $\mathbf{x}_\mu(s)$ and $\mathbf{x}_\sigma(s)$ are feature vectors.

Example: Use an ANN whose output denotes μ and σ

A Revisit to Policy Gradient Theorem

Discrete Action Space

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[\sum_a q_{\pi_{\boldsymbol{\theta}}}(S_t, a) \nabla \pi_{\boldsymbol{\theta}}(a|S_t) \right] \\ &= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} [G_t \nabla \ln \pi_{\boldsymbol{\theta}}(A_t|S_t)]\end{aligned}$$

Continuous Action Space

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[\int_a q_{\pi_{\boldsymbol{\theta}}}(S_t, a) \nabla \pi_{\boldsymbol{\theta}}(a|S_t) da \right] \\ &= \mathbb{E}_{\pi_{\boldsymbol{\theta}}} [G_t \nabla \ln \pi_{\boldsymbol{\theta}}(A_t|S_t)]\end{aligned}$$

We need $\nabla \pi_{\boldsymbol{\theta}}(a|S_t)$!

With the above parametrization with normal pdf, we can easily calculate $\nabla \ln \pi_{\boldsymbol{\theta}}(a|s)$, i.e. $\nabla_{\boldsymbol{\theta}_\mu} \ln \pi_{\boldsymbol{\theta}}(a|s)$ and $\nabla_{\boldsymbol{\theta}_\sigma} \ln \pi_{\boldsymbol{\theta}}(a|s)$.

For instance,

$$\begin{aligned}\nabla_{\boldsymbol{\theta}_\mu} \ln \pi_{\boldsymbol{\theta}}(a|s) &= \nabla_{\boldsymbol{\theta}_\mu} \left(-\frac{(a - \mu(s, \boldsymbol{\theta}))^2}{2\sigma(s, \boldsymbol{\theta})^2} \right) \\ &= \frac{1}{\sigma(s, \boldsymbol{\theta})^2} (a - \mu(s, \boldsymbol{\theta})) \mathbf{x}_\mu(s)\end{aligned}$$

Self-study Exercise: Find $\nabla_{\boldsymbol{\theta}_\sigma} \ln \pi_{\boldsymbol{\theta}}(a|s)$.

Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Recap: Pros and Cons of Policy Gradient Methods

Pros:

- ▶ Deal with high-dimensional or continuous action spaces directly
- ▶ Directly learn stochastic policies
- ▶ Better convergence properties
 - ▶ directly approximate gradient ascent on the performance
- ▶ Incorporate a prior information about the problem

Cons:

- ▶ Typically converge to a local rather than global optimum
- ▶ Typically sample inefficient and high variance

Discussions: Different Policy Gradient Ideas from the Literature -A2C

Advantage Actor Critic (A2C):

Similar to the before:

- ▶ estimate advantage function (but with n-step returns)
- ▶ a shared ANN to parametrize π_{θ} and \hat{v}_{θ}
- ▶ multiple-thread processing
- ▶ add an exploration term to the objective to prevent early convergence to deterministic policies

Discussions: Different Policy Gradient Ideas from the Literature - TRPO

Notation:

$$r_t(\theta) = \frac{\pi_\theta(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}$$

Advantage function: $A_{\theta_{old}} = A_{\theta_{old}}(S_t, A_t)$

Trust Region Policy Optimization (TRPO):

$$\max_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} [r_t(\theta) A_{\theta_{old}}]$$

st.

$$\mathbb{E}_{\pi_{\theta_{old}}} [D_{KL}(\pi_{\theta_{old}}(\cdot|S_t) || \pi_\theta(\cdot|S_t))] \leq \epsilon_{TRPO}$$

Discussions: Different Policy Gradient Ideas from the Literature -PPO

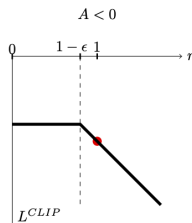
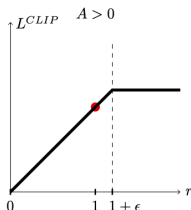
Notation:

$$r_t(\theta) = \frac{\pi_\theta(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}$$

Advantage function: $A_{\theta_{old}} = A_{\theta_{old}}(S_t, A_t)$

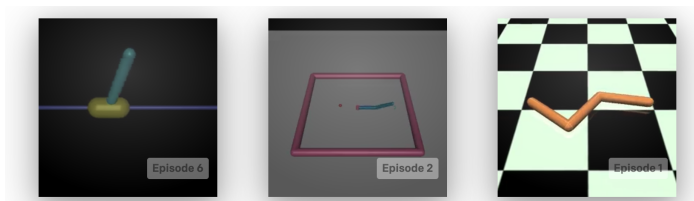
Proximal Policy Optimization (PPO):

$$\max_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} [\min(r_t(\theta)A_{\theta_{old}}, \text{clip}(r_t(\theta), 1 + \epsilon, 1 - \epsilon)A_{\theta_{old}})]$$



IDEA: obtain a lower bound on the objective by ignoring the change in r_t when it would make the objective improve too much, and include it when it makes the objective worse

Mujoco: Continuous control tasks running in a physics simulator

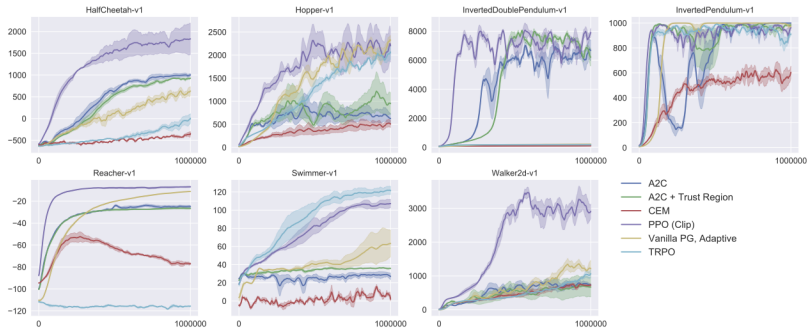


Inverted Pendulum

Reacher

Swimmer

Comparison of Different Policy Gradient Algorithms



Course Map

What are Policy Gradient Methods?

Motivation: Why are policy gradient methods attractive?

Policy Optimization

REINFORCE: Monte-Carlo Policy Gradient

REINFORCE with Baseline

Actor-Critic

Policy Parametrization for Continuous Actions

Discussions

Concluding Remarks

Summary - Key Concepts

Parametrization of the policy: soft-max policies for discrete action spaces, Gaussian parametrization for continuous actions

Aim: Find the best policy by maximizing $J(\theta) = v_{\pi_\theta}(s_0)$ over θ :

$$\max_{\theta \in \mathbb{R}^{d'}} J(\theta)$$

Main Idea for Policy Gradient Methods:

$$\theta_{t+1} = \theta_t + \alpha^\theta \nabla J(\theta_t)$$

Advantage Function:

$$A_{\pi_\theta}(s, a) = q_{\pi_\theta}(s, a) - v_{\pi_\theta}(s)$$

Policy Gradient:

$$\nabla J(\theta) = \mathbb{E}_{\pi_\theta} [q_{\pi_\theta}(S_t, A_t) \nabla \ln \pi_\theta(A_t | S_t)] = \mathbb{E}_{\pi_\theta} [A_{\pi_\theta}(S_t, A_t) \nabla \ln \pi_\theta(A_t | S_t)]$$

Concluding Remarks -Comparison

REINFORCE

$$\theta \leftarrow \theta + \alpha^\theta G_t \nabla \ln \pi_\theta(A_t|S_t)$$

- high variance, updates unbiased, inconvenient for online/continuing problems

REINFORCE with a state-dependent baseline $\hat{v}(S_t, \mathbf{w})$ found using G_t

$$\theta \leftarrow \theta + \alpha^\theta (G_t - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_\theta(A_t|S_t)$$

- lower variance, updates unbiased, inconvenient for online/continuing problems

One-step Actor-Critic

$$\theta \leftarrow \theta + \alpha^\theta (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla \ln \pi_\theta(A_t|S_t)$$

- lower variance, updates biased, convenient for online/continuing problems

References

- ▶ Sutton, Barto, Reinforcement Learning: An Introduction
- ▶ V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. "Asynchronous methods for deep reinforcement learning". PMLR, 2016
- ▶ J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. "Trust region policy optimization", PMLR, 2015
- ▶ J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms", 2017