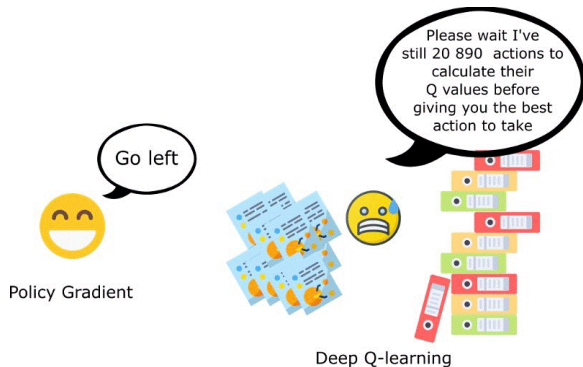


Policy gradient

Wen Yu

Departamento de Control Automático
CINVESTAV-IPN

Problems of value based method



Problems of value based method

- 1 DQN generally only deal with discrete actions and cannot deal with continuous actions.
- 2 Due to the limitations of individual observations or the limitations of modeling, two states that are originally different in the real environment have the same feature description after we model them. Value-based method cannot get the optimal solution.
- 3 The optimal strategy corresponding to the value-based reinforcement learning method is usually a deterministic strategy, while the optimal strategy for some problems is a random strategy. In this case, it is also impossible Solved through value-based learning.
- 4 the policies may not exist with the action-value estimation

Advantages of policy approximation

- Parameterizing policies with the soft-max in action preferences can approach a deterministic policy.
- Action-value based methods have no natural way of finding stochastic optimal policies, whereas policy approximating methods can.
- The soft-max in action preferences enables the selection of actions with arbitrary probabilities.

Policy-Based Reinforcement Learning

The action-value function can be approximated as

$$\begin{aligned}\hat{V}(s, w) &\approx V^\pi(s) \\ \hat{Q}(s, a, w) &\approx Q^\pi(s)\end{aligned}$$

A policy was generated directly from the value function, ϵ -greedy policy

$$\pi(a, s) = \begin{cases} \arg \max_a Q^\pi(s, a) & \text{probability } (1 - \epsilon) \\ a & \text{probability } \frac{\epsilon}{|\mathcal{A}|} \end{cases}$$

Now we will directly parametrize the policy (not table policy)

$$\pi(a, s, w) = P(a | s, w)$$

The goal is to find $\pi(a, s, w)$ such that

$$\max V^\pi(s), \quad \max_{\theta} \mathbb{E}[R | \pi_{\theta}]$$

Policy-Based vs. Value-Based

- Value based: learn value function, implicit policy (e.g, ϵ -greedy)
- Policy base: no value function, learn policy
- Actor-Critic: learn value function, learn policy

Policy search methods

- Policy search methods: model-free and model-based
- Model-free policy: random search and deterministic search
- Stochastic search: developed as the gradient method
- Policy gradient: difficulty in determining the learning rate.

Quality of a policy

The quality of a policy $\pi(\theta)$:

- 1 value of the policy in an episodic

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

- 2 The average value

$$J_{av} = \sum_s P^{\pi(w)}(s) V^{\pi(\theta)}(s)$$

where $P^{\pi(\theta)}$ is distribution of the Markov chain for $\pi(w)$

- 3 The average reward per time step

$$J_{avR} = \sum_s P^{\pi(\theta)}(s) \sum_a \pi(\theta) R(a, s)$$

Approximate the gradient

Episodic MDP

The goal is to find $\pi(a, s, \theta)$ such that $\max V^\pi(s)$ by gradient of the policy, w.r.t, parameter w

$$\Delta\theta = \alpha \nabla_\theta V(\theta)$$

where α is a step.size parameter, $\nabla_w V(w)$ is the policy gradient,

$$\nabla_\theta V(\theta) = \left[\frac{\partial V(\theta)}{\partial \theta_1} \dots \frac{\partial V(\theta + \epsilon)}{\partial \theta_n} \right]$$

we can use

$$\frac{\partial V(\theta)}{\partial \theta_i} \approx \frac{V(\theta_i + \epsilon) - V(\theta_i)}{\epsilon}$$

to approximate the gradient $\nabla_\theta V(\theta)$

Compute the policy gradient analytically

Policy value of π_θ

$$V(\theta) = \mathbb{E}_{\pi_\theta} \left(\sum R_{\pi_\theta} \right) = \sum_i P(i, \theta) R(i)$$

our goal is

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_i P(i, \theta) R(i)$$

Take gradient with respect to θ , using $\frac{\partial}{\partial x} \log x = \frac{1}{x}$

$$\begin{aligned} \nabla_{\theta} V(\theta) &= \nabla_{\theta} \sum_i P(i, \theta) R(i) = \sum_i R(i) \nabla_{\theta} P(i, \theta) \\ &= \sum_i R(i) \frac{P(i, \theta)}{P(i, \theta)} \nabla_{\theta} P(i, \theta) \\ &= \sum_k P(i, \theta) R(i) \nabla_{\theta} \log P(i, \theta) \end{aligned}$$

Approximate

$$\nabla_{\theta} V(\theta)$$

The policy gradient becomes the expectation and we can use the empirical average to estimate it. We use the current policy to sample m trajectories, the empirical average of these m trajectories can be used to approximate the policy gradient:

Approximate $\nabla_{\theta} V(\theta)$ with m sample path under the policy π_{θ} (same $P(i, \theta)$)

$$\frac{\partial V(\theta)}{\partial \theta_i} \approx \frac{1}{m} \sum_k R(i) \nabla_{\theta} \log P(i, \theta) = f(x_i) \nabla \log P(x_i | \theta)$$

so

$$\Delta \theta = \alpha f(x_i) \nabla \log P(x_i | \theta)$$

moving in the direction to yjr logprob of the sample

The item $\nabla \log P(x_i | \theta)$ is the steepest direction in which the probability of the trajectory changes with the parameter θ .

$\frac{1}{m} \sum_k R(i)$ controls the direction and step length of parameter update.

Calculate the steepest direction

Calculate $\nabla \log P(x_i | \theta)$, for i -th episode,

$$P(x_i | \theta) = \prod P(s_{t+1}^i | s_t^i, a_t^i) \pi_\theta(s_t^i | a_t^i)$$

and

$$\log P(x_i | \theta) = \sum \log P(s_{t+1}^i | s_t^i, a_t^i) + \sum \log \pi_\theta(s_t^i | a_t^i)$$

So

$$\begin{aligned} \nabla_\theta \log P(x_i | \theta) &= \sum \nabla_\theta \log P(s_{t+1}^i | s_t^i, a_t^i) + \sum \nabla_\theta \log \pi_\theta(s_t^i | a_t^i) \\ &= \sum \nabla_\theta \log \pi_\theta(s_t^i | a_t^i) \end{aligned}$$

The likelihood gradient $\nabla \log P(x_i | \theta)$ is transformed into the gradient of the action $\nabla_\theta \log \pi_\theta(s_t^i | a_t^i)$,

But it has nothing to do with dynamic model

$$P(s_{t+1}^i | s_t^i, a_t^i)$$

Calculate the scoll function

because

$$\begin{aligned} J_{avR} &= \sum_s P^{\pi(w)}(s) \sum_a \pi(w) R(a, s) \\ \nabla J_{avR} &= \nabla_{\theta} V(\theta) = \frac{\partial V(\theta)}{\partial \theta_i} \\ &\approx \frac{1}{m} \sum_{k=1}^m R(T^i) \sum_{i=1}^T \nabla_{\theta} \log \pi_{\theta}(s_t^i | a_t^i) \end{aligned}$$

so

$$\nabla_{\theta} V(\theta) \approx \frac{1}{m} \sum_k R(i) \sum_k \nabla_{\theta} \log \pi_{\theta}(s, a)$$

We need to calculate $\nabla_{\theta} \log \pi_{\theta}(s, a)$ So

$$\nabla_{\theta} J_{avR} = \mathbb{E} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi}]$$

The term $\frac{1}{m} \sum_{k=1}^m$ likes MC, is unbiased, but the variance is large (noisy).

Temporal structure: Reinforce Algorithm

Like TD learning or Q learning

$$\begin{aligned}\nabla_{\theta} V(\theta, r) &= \mathbb{E} \left[\left(\sum_{i=0}^{T-1} r_i \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(s, a) \right) \right] \\ \nabla_{\theta} V(\theta, r') &= \mathbb{E} \left[\left(\sum_{i=0}^{T-1} r'_i \right) \left(\sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(s, a) \right) \right]\end{aligned}$$

But $\sum_{i=0}^{T-1} r_i = G_t$,

$$\begin{aligned}\nabla_{\theta} V(\theta) &\approx \frac{1}{m} \sum_{k=1}^m R(T^i) \sum_{i=1}^T \nabla_{\theta} \log \pi_{\theta}(s_t^i | a_t^i) \\ &= \frac{1}{m} \sum_{k=1}^m \sum_{i=1}^T \nabla_{\theta} \log \pi_{\theta}(s_t^i | a_t^i) G_t\end{aligned}$$

For $t=1$ to $T-1$

$$\theta(t+1) = \theta(t) + \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) G_t$$

Calculate logPi: Gaussian

The policy can be parameterized in any way, as long as π_θ is differentiable. Generally, a random policy can be written as a deterministic part plus a random part,

$$\pi_\theta = h_\theta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

For the Gaussian distribution ε with a mean value of zero and a standard deviation. The deterministic part is commonly expressed as linear in features

$$h_\theta(s, a) = \theta^T z(s, a)$$

Total is

$$a \sim N(\theta^T z(s, a), \sigma^2)$$

Calculate logPi: example-Soft-max

The action space distribution is soft-max distribution, for each action

$$\pi_{\theta}(a | s) = \frac{e^{z(s,a)}}{\sum_a e^{z(s,a)}}$$

$z(s, a)$ can be parameterized, such as neural network

$$z(s, a) = \phi^T(s, a) \theta$$

then

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

$\phi^T(s, a)$ is current feature, $\mathbb{E}[\phi^T(s, a)]$ is average feature over all actions under the policy

Calculate logPi: example-Gaussian

We fix the variance σ^2 ,

$$\begin{aligned} a &\sim \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{[a-\mu(s)]^2}{2\sigma^2}} \\ \pi_\theta(a | s) &= \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{[a-\mu(s)]^2}{2\sigma^2}} \\ \mu(s) &= \phi^T(s) \theta \\ \nabla_\theta \log \pi_\theta(s, a) &= \frac{[a-\mu(s)]\phi(s)}{\sigma^2} \end{aligned}$$

Reinforce with baseline

Since

$$\nabla_{\theta} V(\theta) \approx \frac{1}{m} \sum_k R(i) \sum_k \nabla_{\theta} \log \pi_{\theta}(s, a)$$

then

$$\Delta \theta_t = \alpha G_t \nabla_{\theta} \log \pi_{\theta}, \quad \text{SGD}$$

The policy gradient is unbiased, but the variance is large. To fix it, we introduce temporal structure or baseline to reduce variance, since

$$V(\theta) = \mathbb{E}_{\pi_{\theta}} (\sum r \pi_{\theta})$$

$$\nabla_{\theta} V(\theta) = \nabla_{\theta} \mathbb{E}_{\pi_{\theta}} (\sum \nabla_{\theta} \log \pi_{\theta}(s, a) (\sum (r - b)))$$

where $b = \mathbb{E}_{\pi_{\theta}} \sum r$, the update law is

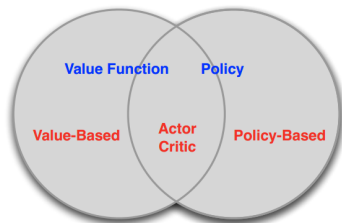
$$\Delta \theta_t = \alpha [G_t - b(s)] \nabla_{\theta} \log \pi_{\theta}, \quad \text{SGD}$$

- Q-learning uses the value of $Q(x, a)$ to take a certain action. It acts as a critic, who evaluates the decision and the evaluation result using $Q_k(x_k, a_k)$
- Q-learning algorithm needs to discretize the action space, which makes the Q-learning algorithm difficult to find the optimal value and the calculation speed is relatively slow.
- Policy gradient calculates the **next action**. The output is the action or distribution of actions. It is an actor.

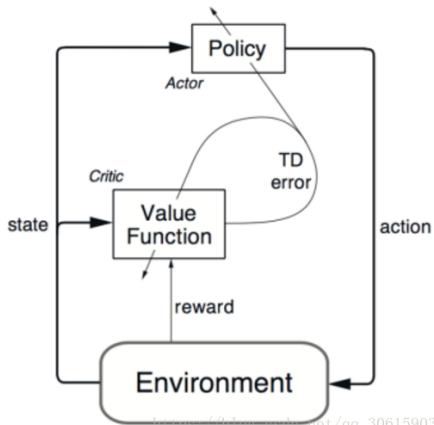
Asynchronous advantage actor-critic

Deep Deterministic Policy Gradient (DDPG): NN->DQN

- Value Based
 - Learnt Value Function
 - Implicit policy (e.g. ϵ -greedy)
- Policy Based
 - No Value Function
 - Learnt Policy
- Actor-Critic
 - Learnt Value Function
 - Learnt Policy



A3C



If the baseline is a critic, then the baseline method becomes actor-critic method.

$$\Delta\theta_t = \alpha [G_t - b(s)] \nabla_{\theta} \log \pi_{\theta}, \quad \text{SGD}$$

The update law is

$$\begin{aligned}\Delta\theta_t &= \alpha [G_t - b(s_t)] \nabla \log \pi_{\theta} \\ &= \alpha [q_t - \hat{V}(s_t, w)] \nabla \log \pi_{\theta} \\ &= \alpha [r_{t+1} + \gamma \hat{V}(s_{t+1}, w) - \hat{V}(s_t, w)] \nabla \log \pi_{\theta}\end{aligned}$$

or

$$\Delta\theta_t = \alpha e_t \nabla \ln [\pi(A_t | s_t, w_t)]$$

where the TD error is

$$e_t = r_{t+1} + \gamma \hat{V}(s_{t+1}, w) - \hat{V}(s_t, w)$$

AC can be regarded as

$$J(w) = \sum_k \ln \pi_w(x_k, u_k) [r + \gamma R_k(x_{k+1}) - R_k(x_k)]$$

Here $[r + \gamma R_k(x_{k+1}) - R_k(x_k)]$ can be regarded as the TD error and it is estimated by the Q-learning algorithm.

Actor-critic: Policy Based+Q-learning

- Learn an approximate value function based.
- Learn approximations to both policy and value functions, actor-critic methods.
 - 1 The 'actor' is a reference to be learned for the policy: Policy Gradients

$$\Delta\theta_t = \alpha [q_t - \hat{V}(s_t, w)] \nabla \log \pi_\theta, \quad \hat{\pi}_w = W_1 \phi(x) \rightarrow NN_1$$

- 2 The 'critic' refers to the learned value function, such as state-value function: Q-learning

$$\hat{V}(s_t, w) = W_2 \phi(x) \rightarrow NN_2$$

TD error based

$$\begin{aligned} e_t &= r + \gamma R_k(x_{k+1}) - R_k(x_k) \\ e_t &= R_{t+1} + \gamma \hat{V}(s_{t+1}, w) - \hat{V}(s_t, w_t) \\ e_t &= R_{t+1} + \gamma \hat{Q}(s_{t+1}, a_{t+1}, w_t) - \hat{Q}(s_t, w_t) \end{aligned}$$

Theorem

The gradient of state-value function can be transform into action-value function

$$\begin{aligned}\Delta\theta_t &= \alpha \nabla J(\theta_t) \\ &= \nabla \sum_k P(i, \theta) R(i) \nabla_{\theta} \log P(i, \theta) \\ &\approx \frac{1}{m} \sum_k R(i) \sum_k \nabla_{\theta} \log \pi_{\theta}(s, a)\end{aligned}$$

Theorem

Policy gradient is

$$\Delta\theta_t = \alpha \nabla J(\theta_t) = \alpha \nabla v_{\pi}(s, \theta)$$

gives an analytic expression for the gradient of performance with respect to the policy parameter, and it does not involve the derivative of the state distribution

$$\begin{aligned}\frac{\partial}{\partial \theta_t} J(\theta_t) &\propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \frac{\partial}{\partial \theta_t} \pi(a | s, \theta_t) \\ \nabla V(\theta) &= \mathbb{E}_{\pi_{\theta}} [q_{\pi}(a, s) \nabla_{\theta} \log \pi(a | s)]\end{aligned}$$

Proof

The object of the proof:

$$J = \mathbb{E} [v \mid \pi_\theta] = v_\pi$$

$$\begin{aligned} \nabla v_\pi(s) &= \nabla (\sum_a \pi(a \mid s) q_\pi(a, s)) \\ &= \sum_a [\nabla \pi(a \mid s) q_\pi(a, s) + \pi(a \mid s) \nabla q_\pi(a, s)] \end{aligned}$$

where $\pi(a \mid s) \nabla q_\pi(a, s)$ is

$$\begin{aligned} \nabla q_\pi(a, s) &= \nabla [\sum_{\bar{s}, r} p(\bar{s}, r \mid s, a) [r + \nabla v_\pi(\bar{s})]] \\ &= \sum_{\bar{s}} p(\bar{s} \mid s, a) \nabla v_\pi(\bar{s}) \end{aligned}$$

where

$$\nabla v_\pi(\bar{s}) = \sum_{\bar{a}} \nabla \pi(\bar{a} \mid \bar{s}) q_\pi(\bar{a} \mid \bar{s}) + \pi(\bar{a} \mid \bar{s}) \sum_{\check{s}} p(\check{s} \mid \bar{s}, \bar{a})$$

where \check{s} is the next step of the state \bar{s} .

$\nabla v_\pi(\bar{s})$ is in regression form,

$$\nabla v_\pi(s) = \sum_a \left[\begin{array}{l} \nabla \pi(a \mid s) q_\pi(a, s) + \\ \pi(a \mid s) \sum_{\bar{s}} p(\bar{s} \mid s, a) \left\{ \begin{array}{l} \sum_{\bar{a}} \nabla v_\pi(\bar{a} \mid \bar{s}) q_\pi(\bar{a} \mid \bar{s}) \\ + \pi(\bar{a} \mid \bar{s}) \sum_{\check{s}} p(\check{s} \mid \bar{s}, \bar{a}) \nabla v_\pi(\check{s}) \end{array} \right\} \end{array} \right]$$

With the parameter θ and a special state s_0

$$\begin{aligned}\nabla J(\theta) &= \nabla v_{\pi}(s_0) = \sum_s \sum_{k \rightarrow \infty} \Pr(s_0 \rightarrow s, k, \pi) \sum_a \nabla \pi(a | s) q_{\pi}(a, s) \\ &= \sum_s \eta(s) \sum_a \nabla \pi(a | s) q_{\pi}(a, s)\end{aligned}$$

where $\sum_s \eta(s)$ can be rewritten as

$$\begin{aligned}\sum_s \eta(s) &= \sum_{\bar{s}} \sum_s \frac{\eta(\bar{s})}{\eta(\bar{s})} \eta(s) = \sum_{\bar{s}} \eta(\bar{s}) \sum_s \frac{\eta(s)}{\eta(\bar{s})} \\ &= \sum_{\bar{s}} \eta(\bar{s}) \sum_s \eta(s)\end{aligned}$$

So

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a \nabla \pi(a | s) q_{\pi}(a, s)$$

It is

$$\frac{\partial}{\partial \theta_t} J(\theta_t) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \frac{\partial}{\partial \theta_t} \pi(a | s, \theta_t)$$