

Policy Evaluation and Value function Approximation

Wen Yu

Departamento de Control Automático
CINVESTAV-IPN

How to learn good policy

- On-policy: evaluate a policy from data obtained from **that** policy. Sampling policy is the same as learning policy. (s_1, a_1, s_2, a_2)
- Off-policy: evaluate a policy from data obtained from a **different** policy. Sampling policy is different with learning policy. $(s_1, a_1, s_2, a_1), (s_1, a_2, s_2, a_2)$

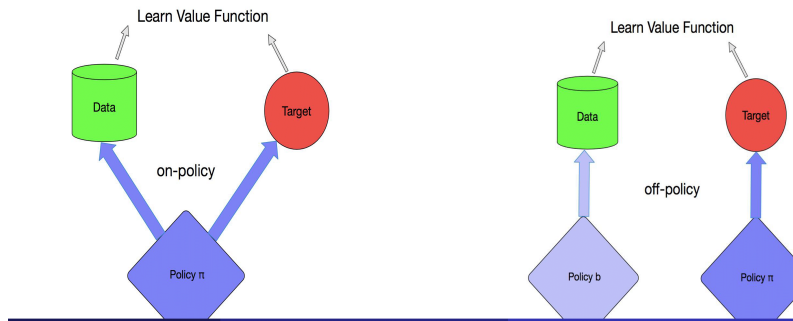
On-policy:

- It learns a better policy using the data from this policy. It is local optimal.
- It cannot assure both exploration and exploitation.

Off-policy approach used two policies:

- One policy generates behavior (exploration, go anywhere). It is called behavior policy, b
- Another policy is for exploitation (use optimal policy). It is called target policy, π

On-policy/Off-policy



Generalized Policy Interaction

Initialize policy π

Repeat

Polity evaluation: compute Q^π

Policy improvement: update π

$$\begin{aligned}\pi'(s) &= \arg \max_a Q^\pi(s, a) \\ &= \arg \max_a [r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^\pi(s')]\end{aligned}$$

It is model-based, we need $p(s' | s, a)$

MC for Q evaluation (On Policy)

MC Policy Evaluation (Model free)

Initialize the counter $N(s, a) = 0$, $G(s, a) = 0$, $Q^\pi(s, a) = 0$, $\forall s \in S$,
 $\forall a \in A$

Loop

Using policy π sample the episode i until T_i , generate

$s_{i,1}, a_{i,1}, r_{i,1}, \dots, s_{i,T_i}, a_{i,T_i}, r_{i,T_i}$,

calculate the return from t to all path of the i th episode

$$G_{i,t} = r_{i,t} + \gamma r_{i,t+1} + \dots + \gamma^{T_i-1} r_{i,T_i}$$

For each state-action (s, a) visited in episode i

For the **first** time t that the state (s, a) is visited in episode i

$$N(s, a) = N(s, a) + 1$$

$$G(s, a) = G(s, a) + G_{i,t}(s, a)$$

$$Q^\pi(s, a) = \frac{G(s, a)}{N(s, a)}$$

Given an estimate $Q^\pi(s, a), \forall s, \forall a$

Update new policy

$$\pi_{t+1}(s) = \arg \max_a Q^{\pi_t}(s, a)$$

$$\arg \max_a Q^{\pi_i}(s, a) \rightarrow a$$

$$\max_a Q^{\pi_i}(s, a) \rightarrow Q$$

Policy Evaluation with Exploration

- Need to try all (s, a) pairs, then follow π
- $Q^\pi(s, a)$ should be good enough so that policy improvement is a monotonic operator

Balance exploration and exploitation. Define $|\mathcal{A}|$ be the number of actions. ϵ -greedy policy is

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

where a is a random action with probability ϵ , but we take a $|\mathcal{A}|$ times, $\max_a Q^\pi(s, a)$ is greedy with probability $(1 - \epsilon)$

It is also called ϵ -soft, or soft policy

Convergence the soft-policy

Theorem

For any ϵ -greedy policy π , w.r.t. Q^π , π_{i+1} is a monotonic improvement

$$V^{\pi_{i+1}} \geq V^{\pi_i}$$

Proof.

Because

$$\begin{aligned} Q^{\pi_k}(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{\pi_k}(s') \\ Q^\pi(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \pi(s' | s, a) Q^\pi(s') \end{aligned}$$

so

$$\begin{aligned} V^{\pi_{i+1}} &= V^\pi(s, \pi_{i+1}) = \mathbb{E}[Q^{\pi_{i+1}}(s, a)] = \sum_{a \in A} \pi_{i+1}(a | s) Q^\pi(s, a) \\ &= (1 - \epsilon) \max_a Q^\pi(s, a) + \frac{\epsilon}{|A|} \sum_{a \in A} Q^\pi(s, a) \end{aligned}$$



Convergence the soft-policy

We calculate

$$\begin{aligned}(1 - \epsilon) \max_a Q^\pi(s, a) &= (1 - \epsilon) \max_a Q^\pi(s, a) \frac{1 - \epsilon}{1 - \epsilon} \\ &= (1 - \epsilon) \max_a Q^\pi(s, a) \frac{(\sum_a \pi_i(a|s)) - \epsilon}{1 - \epsilon} \\ &\geq (1 - \epsilon) Q^\pi(s, a) \frac{\sum_a [\pi_i(a|s) - \frac{1}{|\mathcal{A}|} \epsilon]}{1 - \epsilon} \\ &= Q^\pi(s, a) \sum_a \left[\pi_i(a|s) - \frac{1}{|\mathcal{A}|} \epsilon \right] \\ &= \sum_a \pi_i(a|s) Q^\pi(s, a) - \frac{\epsilon}{|\mathcal{A}|} \sum_a Q^\pi(s, a)\end{aligned}$$

Convergence the soft-policy

So

$$\begin{aligned} V^\pi(s, \pi_{i+1}) &\geq \sum_a \pi_i(a | s) Q^\pi(s, a) - \frac{\epsilon}{|\mathcal{A}|} \sum_a Q^\pi(s, a) + \frac{\epsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} Q^\pi(s, a) \\ &= \sum_a \pi_i(a | s) Q^\pi(s, a) = \mathbb{E}[Q^\pi(s, a)] = V^{\pi_i} \end{aligned}$$

So the policy improvement

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

can assure

$$V^{\pi_{i+1}} \geq V^{\pi_i}, \text{ for all } s \in \mathcal{S}$$

ϵ -greedy policy improvement over policy π , for any $s \in \mathcal{S}$

Convergence the soft-policy

When the ϵ -soft policy π is no longer improved, then

$$V^{\pi_{i+1}} = V^{\pi_i}$$

Proof.

When the policy π is no longer improved,

$$\max_a Q^\pi(s, a) = Q^\pi(s, a)$$

so

$$\begin{aligned} V^\pi(s, \pi_{i+1}) &= \frac{\epsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} Q^\pi(s, a) + (1 - \epsilon) \max_a Q^\pi(s, a) \\ &= \sum_a \pi_i(a | s) Q^\pi(s, a) = V^{\pi_i} = V^{\pi^*} \end{aligned}$$



GLIE (Greedy in Limit of Infinite Exploration)

If all state-action pairs are visited an infinite number of times

$$\lim_{i \rightarrow \infty} N_i(s, a) \rightarrow \infty$$

This means $\epsilon = 0$, the ϵ -greedy policy becomes

$$\pi(a | s) = \arg \max_a Q^\pi(s, a)$$

GLIE

$$\lim_{i \rightarrow \infty} \pi_i(a | s) \rightarrow \arg \max_{\text{with prob 1}} Q(s, a)$$

So the behavior policy $\max Q(s, a)$ converges to greedy policy $\pi_i(a | s)$.

MC for Q function evaluation with ϵ -greedy

Initialize the counter $N(s, a) = 0$, $Q(s, a) = 0$, $\forall s \in S$, $\forall a \in A$, set $\epsilon = 1$, $i = 1$
Create initial ϵ -greedy policy

$$\pi_i = \epsilon\text{-greedy}(Q)$$

Loop

sample the episode i with $\pi_i : s_{i,1}, a_{i,1}, r_{i,1}, \dots, s_{i,T_i}, a_{i,T_i}, r_{i,T_i}$
calculate the return from t to all path of the i th episode

$$G_{i,t} = r_{i,t} + \gamma r_{i,t+1} + \dots + \gamma^{T_i-1} r_{i,T_i}$$

For $t = 1 \dots T$

if first-visit (s, a) in episode i ,

$$N(s, a) = N(s, a) + 1$$

$$Q(s_t, a_t) = \frac{G(s, a)}{N(s, a)}$$
$$= Q(s_t, a_t) + \frac{1}{N(s, a)} [G_{i,t} - Q(s_t, a_t)]$$

$i = i + 1$
(CINVESTAV-IPN)

TD for policy iteration (SARSA)

Initial ϵ -greedy policy π randomly ($\epsilon = 1$), initial state $s_t = s_0$

Sample action from policy, Take $a_t \sim \pi(s_t)$

Observer (r_t, s_{t+1})

Loop

Take action $a_{t+1} \sim \pi(s_{t+1})$

Observer (r_{t+1}, s_{t+2})

Update Q with $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Perform policy improvement, ϵ -greedy(Q)

$t = t + 1$

Convergence properties of SARSA

Theorem

SARSA for finite-state and finite-action MDPs converges to the optimal action-value,

$$Q(s, a) \rightarrow Q^*(s, a)$$

if 1) The policy sequence $\pi_t(a | s)$ satisfied the condition GLIE (All state-action pairs are visited an infinite number of times), 2) The step α_t satisfied the Robbins-Munro sequence

$$\begin{aligned}\sum_{t=1}^{\infty} \alpha_t &= \infty \\ \sum_{t=1}^{\infty} \alpha_t^2 &< \infty\end{aligned}$$

Empirically do not use it.

Robbins-Muno sequence

Check the condition

$$0 < a_t \leq 1, \sum_t a_t = \infty.$$

We usually select constant learning rate

$$a_t = \alpha, \quad 0 < \alpha \leq 1$$

When

$$a_t = \frac{\eta}{1 + \frac{1}{\beta}k}, \quad 0 < \eta \leq 1, \quad \beta \gg 1$$

where η is a constant. Since β is very big, for finite time k , the learning rate $\alpha_k \approx \eta$, is constant.

Robbins-Muno sequence

Because

$$\sum_{k=1}^{\infty} \frac{\eta}{1 + \frac{1}{\beta}k} = \infty$$
$$\sum_{k=1}^{\infty} \left(\frac{\eta}{1 + \frac{1}{\beta}k} \right)^2 = \eta^2 \beta^2 \psi(\beta, 1) - \eta^2 < \infty$$

where $\psi(\beta, 1)$ is the Digamma function, it is bounded.

So Q_k converges to zero w.p.1, and hence, Q_k converges to Q^* with probability one.

Maintain state-action estimates and use bootstrapping, use the value of the best function action

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \left[r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

where SARSA is

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

-greedy exploration

Initial $Q(s, a)$, $\forall s \in S$, initial state $s_t = s_0$

Set π_b to be ϵ -greedy w.r.t. Q

Loop

Take action $a_t \sim \pi_b(s_t)$

Observer (r_t, s_{t+1})

Update Q with (s_t, a_t, r_t, s_{t+1})

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \left[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

Perform policy improvement, ϵ -greedy(Q)

$t = t + 1$

Convergence of Q-learning

Lemma

Consider the stochastic process (ζ, Δ, F) , let P be a sequence of increasing, ζ_0 and Δ_0 are P_0 -measurable, and ζ_k, Δ_k and F_k are P_k -measurable, if Δ_k satisfies

$$\Delta_{k+1} = (1 - \zeta_k) \Delta_k + \zeta_k F_k \quad (1)$$

and

- 1 $0 < \zeta_k \leq 1, \sum_k \zeta_k = \infty, \sum_k \zeta_k^2 < \infty$
- 2 $\|E\{F_k | P_k\}\| \leq \kappa \|\Delta_k\| + c_k, \kappa \in (0, 1],$ and c_k converges to zero
- 3 $\text{var}\{F_k | P_k\} \leq K(1 + \kappa \|\Delta_k\|)^2, K$ is a positive constant, $\|\cdot\|$ denotes the maximum norm.

then Δ_k converges to zero with probability one.

Theorem

For the finite MDP (X, U, f, ρ) , the Q-learning algorithm

$$Q_{k+1}(x_k, u_k) = Q_k(x_k, u_k) + \alpha_k \{ \rho(x_k, u_k) + \gamma \min_{u^*} [Q_k(x_{k+1}, u^*)] - Q_k(x_k, u_k) \}$$

converges to the optimal value function Q^* almost surely, if

$$\sum_k \alpha_k = \infty, \quad \sum_k \alpha_k^2 < \infty \quad (2)$$

Convergence of Q-learning

The Q-learning algorithm is

$$\begin{aligned} Q_{k+1}(x_k, u_k) &= Q_k(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \min_{u^*} Q_k(x_{k+1}, u^*) - Q_k(x_k, u_k)] \\ &= (1 - \alpha_k) Q_k(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \min_{u^*} Q_k(x_{k+1}, u^*)] \end{aligned} \quad (3)$$

where $r_{k+1} = \rho(x_k, u_k)$. Define $\Delta_k(x_k, u_k) = Q_k(x_k, u_k) - Q^*$,

$$\Delta_{k+1}(x_k, u_k) = (1 - \alpha_k) \Delta_k(x_k, u_k) + \alpha_k \left(r_{k+1} + \gamma \min_{u^*} Q_k(x_{k+1}, u^*) - Q^* \right)$$

Convergence of Q-learning

Define

$$F_k(x_k, u_k) = r_{k+1} + \gamma \min_{u^*} Q_k(x_{k+1}, u^*) - Q^*$$

Check the condition

$$\|E\{F_k|P_k\}\| \leq \kappa\|\Delta_k\| + c_k, \kappa \in (0, 1)$$

Use the value iteration mapping \mathcal{H} for P_k ,

$$E\{F_k(x_k, u_k)|P_k\} = \mathcal{H}[Q_k(x_k, u_k)] - Q^* = \mathcal{H}[Q_k(x_k, u_k)] - \mathcal{H}(Q^*)$$

Since \mathcal{H} is a contraction, from Lemma

$$\mathcal{H}[Q_k(x_k, u_k)] - \mathcal{H}(Q^*) \leq \gamma\|Q_k(x_k, u_k) - Q^*(x_k, u_k)\| = \gamma\|\Delta_k(x_k, u_k)\|$$

So

$$\|E\{F_k(x_k, u_k)|P_k\}\| \leq \gamma\|\Delta_k(x_k, u_k)\|$$

Convergence of Q-learning

Check the condition

$$\text{var}\{F_k|P_k\} \leq K(1 + \kappa\|\Delta_k\|)^2$$

$$\text{var}\{F_k(x_k, u_k)|P_k\} = \text{var}\left[\left(r_{k+1} + \gamma \min_{u^*} Q_k(x_{k+1}, u^*) - Q^*\right)^2\right]$$

Because r_{k+1} is bounded, and $\Delta_k(x_k, u_k) = Q_k(x_k, u_k) - Q^*$,

$$\text{var}\{F_k(x_k, u_k)|P_k\} \leq K(1 + \gamma\|\Delta_k(x_k, u_k)\|)^2$$

where K is a constant.

Convergence of Q-learning

Check the condition

$$0 < \zeta_k \leq 1, \sum_k \zeta_k = \infty.$$

For Q-learning we usually select constant learning rate

$$\zeta_k = \alpha, \quad 0 < \alpha \leq 1$$

But if

$$\alpha_k = \frac{\eta}{1 + \frac{1}{\beta}k}, \quad 0 < \eta \leq 1, \quad \beta \gg 1$$

where η is a constant. Since β is very big, for finite time k , the learning rate $\alpha_k \approx \eta$, is constant.

Convergence of Q-learning

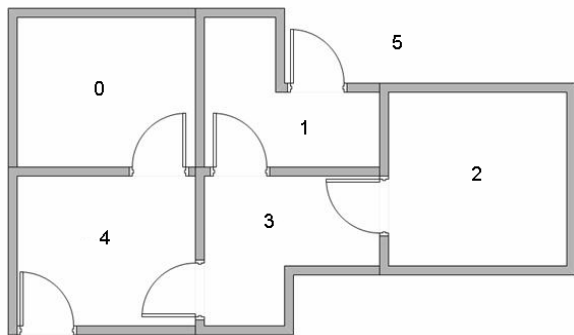
Because

$$\sum_{k=1}^{\infty} \frac{\eta}{1 + \frac{1}{\beta}k} = \infty$$
$$\sum_{k=1}^{\infty} \left(\frac{\eta}{1 + \frac{1}{\beta}k} \right)^2 = \eta^2 \beta^2 \psi(\beta, 1) - \eta^2 < \infty$$

where $\psi(\beta, 1)$ is the Digamma function, it is bounded.

So Δ_k converges to zero w.p.1, and hence, Q_k converges to Q^* with probability one.

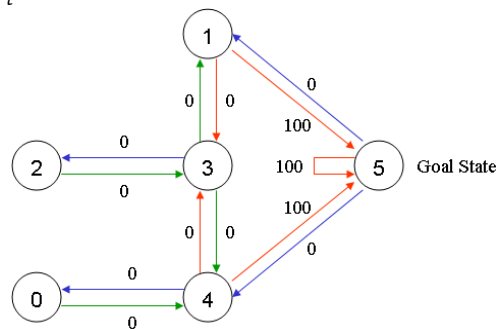
Example



- There are 5 rooms connected by doors.
- The outside of the building is number 5.
- An agent start from any room, and go outside (to Room-5)

Reward R

t



State	Action					
	0	1	2	3	4	5
0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

"-1" represents no link between nodes

Q learning

$$Q_{k+1}(s_t, a_t) = Q_k(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_{a_t} [Q_k(s_{t+1}, a_t)] - Q_k(s_t, a_t) \right]$$

- Each exploration is an episode
- Each episode consists of the agent moving from the initial state to the goal state

The initial state is Room-1, $s_0 = 1$

The initialize matrix is $Q_1(s_0, a) = 0$

let $\gamma = 0.8$, $\alpha = 1$,

$$Q_{k+1}(s_t, a_t) = Q_k(s_t, a_t) + r_{t+1} + \gamma \max_{a_t} [Q_k(s_{t+1}, a_t)] - Q_k(s_t, a_t) \\ = r_{t+1} + \gamma \max_{a_t} [Q_k(s_{t+1}, a_t)]$$

Let $Q(1, 5)$ is the value of Q-table as $Q_{1,5}$. It is means: Room 1 (s_t) \rightarrow Room 5 (s_{t+1}).

$Q_k(s_{t+1}, a)$: from Room 5 there are three routes, we use greedy policy $\max_a \{q(5, 1), q(5, 4), q(5, 5)\}$

$$Q(1, 5) = r(1, 5) + 0.8 \max_a \{Q(5, 1), Q(5, 4), Q(5, 5)\} = 100 + 0.8 * 0 =$$

After one episode

$$Q = \begin{array}{c} \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} \end{array} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Q-learning

Next episode:

The initial state is Room-3, $s_0 = 3$

There are three possible actions: $3 \rightarrow 1$, $3 \rightarrow 2$, $3 \rightarrow 4$

For $3 \rightarrow 1$, there are 2 possible actions: $1 \rightarrow 3$, $1 \rightarrow 5$.

$$Q(3, 1) = r(3, 1) + 0.8 \max_a \{Q(1, 3), Q(1, 5)\} = 0 + 0.8 * 100 = 80$$

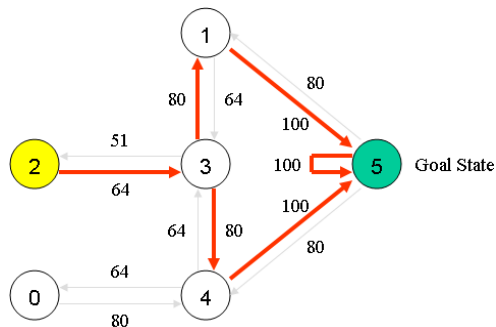
$$Q = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

Then from Room-1, got to Room-5 to finish this episode

After many episodes, Q-table reaches

$$Q = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix} \end{matrix}$$

$$\pi(s_t) = \arg \max_a \{Q_t\}$$



If the initial is Room-2, the optimal trajectory is $2 \rightarrow 3 \rightarrow 1 \rightarrow 5$.

Value function Approximation (VFA)

Generalization: Tabular representation is insufficient

Update the estimator after each episod (MC) or each step (TD): teacher

Value function Approximation (VFA)

Approximating v_π from experience using a known policy π .
Approximate state value,

$$\hat{V}(s, w) \approx V^\pi(s)$$

Approximate control methods: approximation of the action-state value of policy π ,

$$\hat{Q}(s, a, w) \approx Q^\pi(s, a)$$

There are many ways to approximate the value function. For example, in the simplest linear representation, using $\phi(s)$ to represent the eigenvector of the state s , the state value function can be approximately expressed as:

$$\hat{V}(s, w) = W^T \phi(s)$$

- 1 Linear feature representation
- 2 Neural Networks
- 3 Decision trees: highly interpretable
- 4 Nearest neighbors
- 5 Fourier/Wavelet bases

Differential function approximators (smooth optimization property)

- Linear feature representation
- Neural Networks

Linear function approximation

We use a feature vector to represent the feature of the state s ,

$$X(s) = [x_1(s) \cdots x_n(s)]$$

where the features are basis function, which can be defined in many different ways. Value function

$$\hat{V}(s, w) = \sum_{i=1}^n W_i X_i(s) = X^T(s) W$$

Linear function approximation

The object function is

$$J_w = \mathbb{E}_\pi \left\{ [V^\pi(s) - \hat{V}(s, w)]^2 \right\}$$

The real value function $V^\pi(s)$ uses TD (SARSA) or Q-learning.

The stochastic gradient descent is

$$\begin{aligned} \Delta W &= -\alpha \nabla J(W) = -\alpha \frac{\partial J}{\partial W} \\ &= -2\alpha [\hat{V}(s, w) - V^\pi(s)] \frac{\partial}{\partial W} \hat{V}(s, w) \\ &= -2\alpha [\hat{V}(s, w) - V^\pi(s)] X(s) \end{aligned}$$

where α is learning rate, $\hat{V}(s, w) - V^\pi(s)$ is prediction error, $X(s)$ is feature value

From tabular learning to function approximation

Because we do not have real $V^\pi(s)$ as the target in supervised learning, we use estimated values from MC or TD to replace $V^\pi(s)$

MC: the state value is the expected value of the return, we can substitute $V^\pi(s)$ the the return

$$G_t = \sum_i \gamma^{i-1} R_{i,t}$$

after each episod

$$W_{t+1} - W_t = \alpha \left[G_t - X^T(s_t) W_t \right] X(s_t)$$

Bath MC for value function approximation

Let $G(s_i)$ be an unbiased sample of true expected return $V^\pi(s_i)$

$$J = \arg \min_W \sum_{i=1}^N [G(s_i) - W_i X_i(s)]$$

Least square

$$W = (X^T X)^{-1} XG$$

where X is a matrix of the features of each of N states $X_i(s)$

No any Markov assumption

TD learning for value function approximation

TD learning

$$V_{t+1}^{\pi}(s) = V_t^{\pi}(s) + \alpha [\{r_t + \gamma V_t^{\pi}(s_{t+1})\} - V_t^{\pi}(s)]$$

or

$$V^{\pi}(s) = V^{\pi}(s) + \alpha [r + \gamma V^{\pi}(s') - V^{\pi}(s)]$$

The target is

$$r + \gamma V^{\pi}(s')$$

find weights to minimize the mean squared error return $V^{\pi}(s_i)$

$$J_w = \mathbb{E}_{\pi} \left\{ [r + \gamma \hat{V}^{\pi}(s', W) - \hat{V}^{\pi}(s, W)]^2 \right\}$$

Linear TD(0)

$$\begin{aligned}\Delta W &= \alpha [r_t + \gamma \hat{V}^\pi(s', W) - \hat{V}^\pi(s, W)] X(s) \\ &= \alpha [r + \gamma X^T(s') W - X^T(s) W] X(s)\end{aligned}$$

TD(0) linear value function approximation: Policy Evaluation

Initial $t = 1$, $W = \text{rand}$

Loop

Sample (s_t, a_t, r_t, s_{t+1})

Update weights

$$W = W + \alpha \left[r + \gamma X^T(s') W - X^T(s) W \right] X(s)$$

$t = t + 1$

Action-Value approximation

$$\hat{Q}(s, a, W) \approx Q^\pi(s, a)$$

$$J_W = \mathbb{E}_\pi \left\{ [Q^\pi(s, a) - \hat{Q}(s, a, W)]^2 \right\}$$

Use stochastic gradient descent to find a local minimum

$$\Delta W = -\alpha \frac{\partial J}{\partial W}$$

$$X(s, a) = [x_1(s, a) \cdots x_n(s, a)]$$

$$\hat{Q}(s, a, W) = X^T(s, a) W$$

Action-Value approximation

MC

$$W_{t+1} - W_t = \alpha \left[G_t - X^T (s_t, a_t) W_t \right] X (s_t, a_t)$$

SARSA

$$W = W + \alpha \left[r + \gamma X^T (s', a') W - X^T (s, a) W \right] X (s, a)$$

Q-learning

$$W = W + \alpha \left[r + \gamma \max_{a'} \hat{Q} (s', a', W) - \hat{Q} (s, a, W) \right] \nabla_W \hat{Q} (s, a, W)$$

or

$$W = W + \alpha \left[r + \gamma \max \left(X^T (s', a') W \right) - X^T (s, a) W \right] X (s, a)$$

Deep Neural Networks to approximate

- Value function
- Policy
- Model

Neural network

$$\hat{V} = Wz = W\phi(s, a)$$

$$\hat{Q} = Wz = W\phi[V(s, a) \dots]$$

Stochastic gradient decent (SGD)

$$\Delta W = -\alpha \nabla_W J(W)$$

MC method

$$\Delta W = \alpha [G_t - \hat{Q}(s, a, W)] \nabla_W \hat{Q}(s, a, W)$$

SARSAR

$$\Delta W = \alpha [r + \gamma \hat{Q}(s', a', W) - \hat{Q}(s, a, W)] \nabla_W \hat{Q}(s, a, W)$$

Q-learning

$$\Delta W = \alpha \left[r + \gamma \max_a \hat{Q}(s', a, W) - \hat{Q}(s, a, W) \right] \nabla_W \hat{Q}(s, a, W)$$

- Deep Q-Network (DQN) that is the first deep reinforcement learning method proposed by DeepMind. Paper was published on Nature in 2015
- NN is CNN
- Naive DQN has 3 convolutional layers and 2 fully connected layers to estimate Q values directly from images.
- DNN is easily overfitting in online reinforcement learning.