

Lecture 3 Dynamic Programming

What is Dynamic Programming:

① Divide and Conquer + Memory / swap Time strategy
 ② space

How to use DP:

Problems that satisfied:

- Has optimal substructure \Rightarrow ①
- Overlapping sub-problems. \Rightarrow ②

\Rightarrow MDP satisfies both requirements.

① Bellman equation has recursive structure

② Value function is the expectation of G_t , and G_t will be reused.

Policy Iteration:

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$
 $\xrightarrow{\text{set to 0}} \xrightarrow{\text{set to random choose action}} \pi(s) = \frac{1}{|\mathcal{A}|}$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')] \xrightarrow{\text{Bellman Expectation Equation}}$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|) \xrightarrow{\text{end condition: max change < threshold}}$$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)
 $\xrightarrow{\text{improvement}}$

3. Policy Improvement

$\text{policy-stable} \leftarrow \text{true}$

For each $s \in S$:

greedy strategy

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')] \xrightarrow{q(s, \frac{\pi(s)}{a})}$$

If $\text{old-action} \neq \pi(s)$, then $\text{policy-stable} \leftarrow \text{false}$

If policy-stable , then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

↓
goal of policy iteration is to find
the optimal value-table (V^*) and optimal
action table / policy (π^*)

why the algorithms above is correct?

① Policy Evaluation:

The algorithm is mean to recursively calculate $V(s)$ using Bellman Expectation Equation. As the BEE can estimate the value function $V(s)$, thus it's correct. When it converges, it can estimate the exact value function, but it takes a lot of value iteration. That's why we introduce the δ to determine the accuracy of estimation.

② Policy Improvement:

The greedy strategy is mean to perform the Bellman Optimality Equation: $q_{\pi}^*(s, a) = \max_{a' \in A} q_{\pi}(s, a')$.

For the fact that we have: $V_{\pi}^*(s) = q_{\pi}^*(s, a)$, and for $V_{\pi}^*(s)$, π is the optimal policy. Therefore, this algorithm can improve policy π to optimal policy π^* .

Problem:

Policy iteration updates policy for each Policy Evaluation

Step. It may take many iterations to break the loop.

That makes it relatively slow. \Rightarrow

Improvement 1: introduce δ

Improvement 2: Value Iteration (update policy every iteration)

Value Iteration:

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
 Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop: \Rightarrow Iteration loop

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Get the optimal value function $V^(s)$,
 retrieve its policy and it's the optimal policy*

Why the algorithms above is correct?

The algorithm is meant to calculate the optimal value function $V^*(s)$ through Bellman Optimality Equation.

As the fact that $V^*(s)$ is policy is the optimal policy π^* , thus

$$\pi^* = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma \cdot V^*(s')]$$

Synchronous Dynamic Programming Algorithms

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

- Algorithms are based on state-value function $v_\pi(s)$ or $v_*(s)$
- Complexity $O(mn^2)$ per iteration, for m actions and n states
- Could also apply to action-value function $q_\pi(s, a)$ or $q_*(s, a)$
- Complexity $O(m^2n^2)$ per iteration

$\#V(s)$ BEE
 $O(n)$ $O(mn)$
 $\Rightarrow O(mn) \cdot O(mn) \Rightarrow BEE$, 通过二层的分叉树理解
 $\#q_{\pi}(s, a)$

Advantages and Drawbacks of DP:

Advantages :

- Fully solves the problem
- Effective for medium-sized (几百万当量的 states) problems.

Drawbacks

- Need to have advanced knowledge of MDP (R, P)
- Curse of dimensionality (can't have enough memory)
in large problems.



I. Use sample backups method (e.g. MC, TD)

to solve \Rightarrow

- Model-free (no knowledge of MDP)
- Constant cost of backup (break the curse)

II. Use approximate DP

Lecture 4 MC, TD for prediction

Monte Carlo

First-visit MC prediction, for estimating $V \approx v_\pi$

Initialize:

$\pi \leftarrow$ policy to be evaluated

$V \leftarrow$ an arbitrary state-value function

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Repeat forever: $\xrightarrow{\text{from start state } s_0 \text{ to terminal state } s_t}$

Generate an episode using π

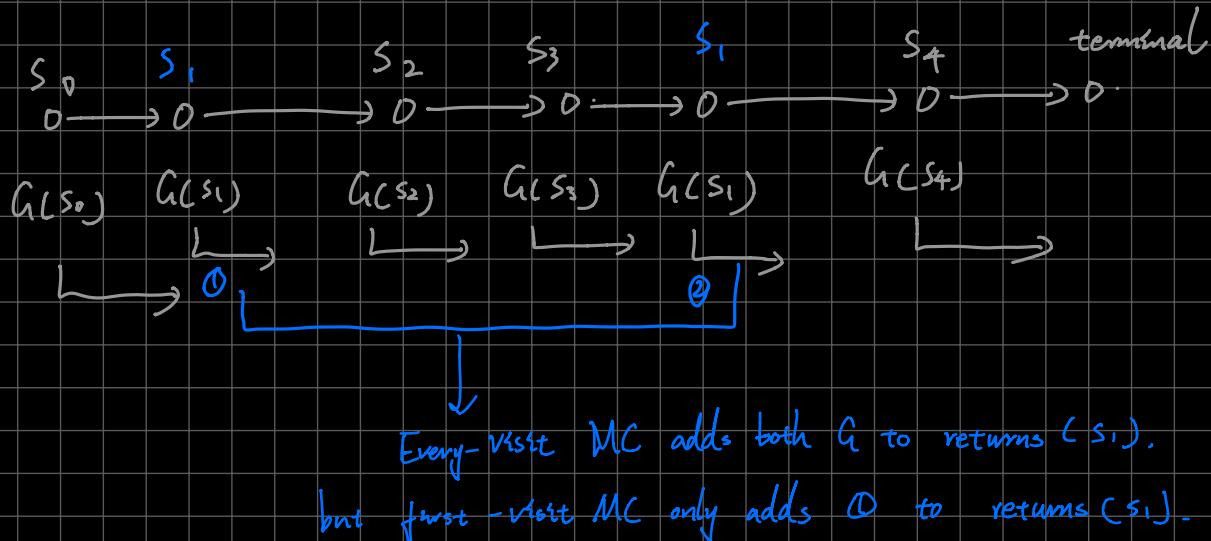
For each state s appearing in the episode:

$G \leftarrow$ return following the first occurrence of $s \Rightarrow$ every visit MC append $G(s)$ to $Returns(s)$

Append G to $Returns(s)$

$V(s) \leftarrow \text{average}(Returns(s))$

Actually this step can be performed outside the loop. 但由于
储存 Returns(s) 的空间太大，因此对每个 s 只执行一次
储存每次优化过的 Returns(s)



Note: MC use simplest possible idea:

Estimated value = empirical average sample returns

使用经验平均(empirical average)估计期望值。
采样 sampling

$$V(s) = E[G(s)]$$

为 value function 定义。

Increment MC:

Original:

节省储存空间，
提高运算精确度

$$\left\{ \begin{array}{l} N(s_t) \leftarrow N(s_t) + 1 \\ V(s_t) \leftarrow \frac{1}{N(s_t)} \cdot \underbrace{\left[G_t + V(s_t) \cdot N_{old}(s_t) \right]}_{\text{sum of Returns}(s_t)} = V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t)) \end{array} \right.$$

Incremental: 这就是增量更新
公式的来源

$$V(s_t) \leftarrow V(s_t) + \underbrace{\alpha (G_t - V(s_t))}_{\text{constant value}}$$

Temporal Difference: (TD(0))

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

$A \leftarrow$ action given by π for S

 Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

 until S is terminal

Note: TD use the bootstrapping method to learn from incomplete episodes. This makes it more efficient and more general (can deploy in non episodic problems)

Difference between MC and TD:

MC: updates value toward actual return G_t

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t)) \quad \text{MC error}$$

TD: updates value toward estimated return: $R_{t+1} + r \cdot V(s_{t+1})$

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + r \cdot V(s_{t+1}) - V(s_t))$$

TD error

$$= \text{TD target} - \text{TD predict}$$

Advantages and Disadvantages:

① TD can learn without the final outcome \Rightarrow
more efficient and more general

② TD has biased, low variance estimate

MC has unbiased, high variance estimate

\Rightarrow MC is not very sensitive to initial values,
but TD is sensitive.

③ TD exploits Markov property while MC doesn't. \Rightarrow

{ TD is more efficient in MC envs,

{ MC is more efficient in non-MC envs,

