

RL Notes - Module 3 Learning Objectives

- Specifying Policies

→ Stochastic vs Deterministic policies

- Policy is a distribution over actions for each possible state

Deterministic policy notation

$$\pi(s) = a$$

States $\xrightarrow{\text{map}}$ Actions

→ for taking each action

```

    graph LR
      s0 --> a0
      s1 --> a1
      s2 --> a2
  
```

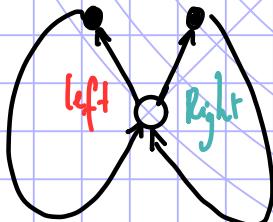
only depend on current state

Stochastic policy notation

$$\sum_{a \in A(s)} \pi(a|s) = 1, \quad \pi(a|s) \geq 0$$

→ Generate valid policies for a given MDP

- Valid and invalid policies



Value functions

- In RL, reward captures the notion of short-term gain.
- The objective is to learn a policy to achieve the most reward on the long run.

⇒ Value funcs formalize what this means.

- Relationship between value funcs and policies

- Brief story of RL (Interesting)

- Bellman Equation Derivation

→ to formalize between value of state and its possible successors

$$V_n(s) \stackrel{def}{=} E_n[G_t | S_t = s], G_t = \sum_{k=0}^{\infty} \gamma^k P_{t+k+1}$$

$$= E_n[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma V_n(s')]$$

- State-value and action-value funcs under policy

$$V_n(s) \stackrel{def}{=} E_n[G_t | S_t = s]$$

↓
 ↗ not incentive
 ↗ random

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$q_n(s, a) \stackrel{def}{=} E_n[G_t | S_t = s, A_t = a]$$

←
 after taking action

$$\rho(\text{win}) = V_n(s) \rightarrow \text{how to be computed?}$$

- Valid value functions

$$E_n[G_{t+1} | S_{t+1} = s']$$

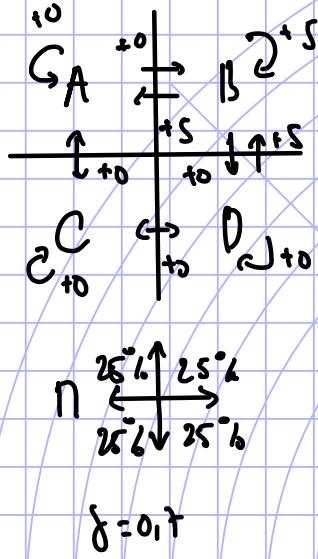
↑
 $E_n[G_t | S_t = s]$

$$q_n(s, a) \stackrel{def}{=} E_n[G_t | S_t = s, A_t = a]$$

(a is fixed)

$$= \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \sum_a \pi(a | s') \cdot q_n(s', a)]$$

E.g. Gridworld



$$\gamma = 0.7$$

$$V_{\pi}(s) = \sum_a n(a|s) \sum_{\pi} \sum_{s'} p(s', r | s, a) \cdot [r + \gamma V_{\pi}(s')]$$

for each a , there's only one possible associated next s' and r

$$= \sum_a n(a|s) \cdot [\pi \cdot \gamma V_{\pi}(s')]$$

$$\left\{ \begin{array}{l} V_{\pi}(A) = \frac{1}{4} (S + 0.7 V_{\pi}(B)) + \frac{1}{4} \cdot 0.7 V_{\pi}(C) + \frac{1}{4} 0.7 \cdot V_{\pi}(D) \\ \dots \end{array} \right. \Rightarrow \left\{ \begin{array}{l} V_{\pi}(A) = 4.2 \\ V_{\pi}(B) = 4.1 \\ V_{\pi}(C) = 2.2 \\ V_{\pi}(D) = 4.2 \end{array} \right.$$

4 linear equations

\Rightarrow Large MDP states \Rightarrow won't be practical

E.g. Chess 10^{45} equations?

- Optimal policies

→ Define an optimal policy

An optimal policy π_* is as good as or better than all the other policies

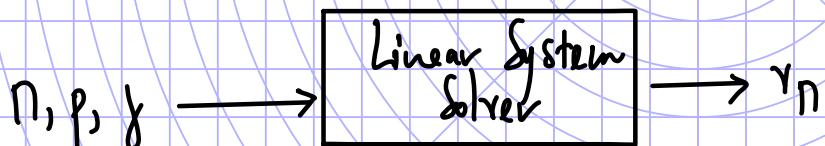
- Optimal value functions

Recall that $\pi_1 > \pi_2 \Leftrightarrow v_{\pi_1}(s) \geq v_{\pi_2}(s)$
for all $s \in S$

$$v_{\pi_*}(s) = E_{\pi_*}[G_t | s_t = s] = \max_{\pi} v_{\pi}(s)$$

for all $s \in S$

$$q_{\pi_*} q_{\pi_*}(s, a) = \max_{\pi} q_{\pi}(s, a) \text{ for all } s \in S \text{ and } a \in A$$



→ Bellman Optimal Equation

$$v_*(s) = \max_a \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_*(s')]$$

$$q_*(s, a) = \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

Using Optimal Value functions to Get Optimal Policies

- The connection between the optimal value function and optimal policies
- Verify off

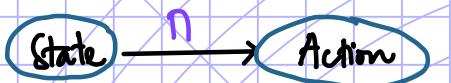
$$\pi_*(s) = \arg \max_a \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_*(s')]$$

$$\pi_*(s) = \arg \max_a q_*(s, a)$$

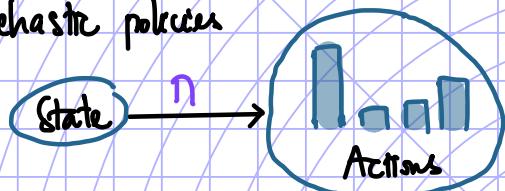
→ Summary

- Policies tell an agent how to behave in their environment

- Deterministic policies



- Stochastic policies



A policy depends only on the **current state**

- The goal: to find a policy that obtains as much reward as possible

An **optimal policy** achieves the highest value possible in every state

$$v_*$$
 $v_{\pi_*}(s) = \max_{\pi} v_{\pi}(s) \text{ for all } s \in \mathcal{S}$

$$= \max_a \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_*(s')]$$

$$q_*$$
 $q_{\pi_*}(s, a) = \max_{\pi} q_{\pi}(s, a) \text{ for all } s \in \mathcal{S} \text{ and } a \in \mathcal{A}$

$$= \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

- Value functions estimate future return under a specific policy

$$v_{\pi}(s) \doteq E_{\pi}[G_t | S_t = s]$$

$$q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a]$$

- Bellman equations define a relationship between the value of a state, or state-action pair, and its possible successors

$$v_{\pi}(s) \doteq \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$q_{\pi}(s, a) \doteq \sum_{s'} \sum_r p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') \cdot q_{\pi}(s', a') \right]$$

$$v_*(s) = \max_a q_*(s, a)$$

Module 4

• Policy Evaluation vs. Control

(1) (2)

- (1) is the task of determining the value function for a specific policy. \rightarrow necessary first step
- (2) is the task of finding a policy to obtain as much reward as possible.
 \hookrightarrow ultimate goal of RL

• Dynamic programming applied, limitations

• Iterative Policy Evaluation

\rightarrow PP are obtained by turning Bellman equations into update rules.

$$v_{k+1}(s) \leftarrow \sum_a n(a|s) \sum_{s'} \sum_r p(s'|r|s,a) [r + \gamma v_k(s')]$$

sweep until $v_k = v_{k+1}$

for any v_0 , $\lim_{k \rightarrow \infty} v_k = v_\pi$

(1) $\pi \rightarrow v_\pi$

Linear System Solver \rightarrow Dynamic Programming

(2) improve a policy \rightarrow strictly better

$\pi, \gamma \rightarrow DP \rightarrow v_\pi$ Policy Evaluation

$\pi, \gamma \rightarrow DP \rightarrow \pi_\pi$ Control

\hookrightarrow Along with Knowledge of π

• Iterative Policy Evaluation, for estimating $v \approx v_\pi$

Input π , the policy to be evaluated

$$v \leftarrow \vec{0}, v' \leftarrow \vec{0}$$

loop:

$$\Delta \leftarrow 0$$

loop for $s \in \mathcal{S}$:

$$v'(s) \leftarrow \sum_a n(a|s),$$

\sum

$$s', r p(s', r | s, a) [r + \gamma v(s')]$$

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

$$v \leftarrow v'$$

until $\Delta < \theta$ (a small positive number)

Output $v \approx v_\pi$

• Policy Improvement

$$\pi_*(s) = \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$

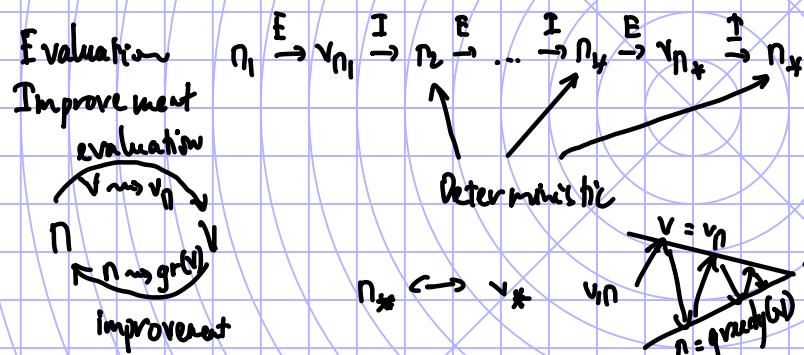
• Policy Improvement Theorem

$$\pi' > \pi: q_\pi(s, \pi'(s)) \geq q_\pi(s, \pi(s)) \text{ for all } s \in S$$

$$\pi' > \pi: q_\pi(s, \pi'(s)) > q_\pi(s, \pi(s)) \text{ for at least one } s \in S$$

• Policy Iteration

→ the dance of policy and value



• Generalized Policy Iteration unified classification

→ Flexibility of the Policy Iteration framework

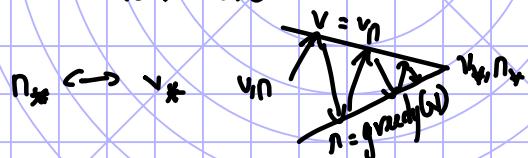
Synchronous vs Asynchronous

→ Value Iteration

$$\text{iff: } V_s \leftarrow \max_a \sum_{s', r} p(s', r | s, a), [r + \gamma V(s')] \text{ and } V(\text{terminal}) = v$$

→ Avoid full Sweeps

→ Update the states near those that have recently changed value



• Efficiency of Dynamic Programming

→ learning a value function:

Monte-Carlo Sampling → Average

Recall

$$v_n(s) = E_n[G_t | S_t = s]$$

→ finding an optimal policy: Brute Force Search

$|A|^{|S|}$ policies

$$v_{k+1}(s) \leftarrow \sum_a n(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_k(s')]$$

• Advantage of DP and bootstrapping over these other strategies

- much faster

→ RL to solve high-dim problem

• ADP for fleet optimization

$$\bar{v}^n(a) = (1-\alpha) \bar{v}^{n-1}(a) + \alpha \hat{v}^n(a)$$

attribute of car

new est old est marginal value

ADP for fleet optimization

Step 1: Start with a pre-decision state S_t^n

Step 2: Solve the deterministic optimization using an approximate value function:

$$\hat{v}_t^n = \min_x (C_t(S_t^n, x_t) + \bar{v}_{t-1}^{n-1}(S_{t-1}^{M,x}(S_t^n, x_t)))$$

to obtain x^t .

Deterministic optimization

Step 3: Update the value function approximation

$$\bar{v}_{t-1}^n(S_{t-1}^{M,x}) = (1 - \alpha_{n-1}) \bar{v}_{t-1}^{n-1}(S_{t-1}^{M,x}) + \alpha_{n-1} \hat{v}_t^n$$

Recursive statistics

Step 4: Obtain Monte Carlo sample of W_t and compute the next pre-decision state:

$$S_t^n = S^M(S_t^n, x_t^n, W_{t+1}(\omega^n))$$

Simulation

Step 5: Return to step 1.

$$\text{merge old with new} \rightarrow \bar{v}^2(\tau x) = (1-\alpha) \bar{v}^1(\tau x) + \alpha \hat{v}^2(\tau x)$$

lets take $\alpha = v_t / t$

$$|\text{States}| = \binom{\text{No. tracks} + |\text{Locations}| - 1}{|\text{Locations}| - 1}$$

$$x_t = \begin{bmatrix} \text{Time} \\ \text{Location} \\ \text{Bicycle} \\ \text{Hours} \end{bmatrix}$$

Linear Programming

→ Cal marginal value $\bar{v}(a)$

→ Multiattribute resources

$$V^+(s_t) = \max \left(C_t(s_t, x_t) + \sum_{\ell \in A} \bar{v}_{t,\ell} \cdot R_{t,\ell} \right)$$

→ hierarchical learning

aggregate vs. disaggregate vs. weighted

Warren Power

$$\begin{bmatrix} a_1 \\ a_2 \\ \dots \\ a_n \end{bmatrix} \longrightarrow w_a^0 \vec{v}^0 G_{11} + w_a^1 \vec{v}^1 \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} + \dots$$

An Approx DP Algo for

Large-Scale Fleet Management : A case application

→ Summary

Policy evaluation is the task of determining the state-value func v_π , for a particular policy π

Iterative Policy evaluation

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_\pi(s')]$$

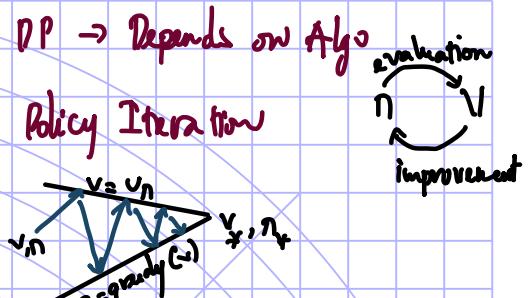
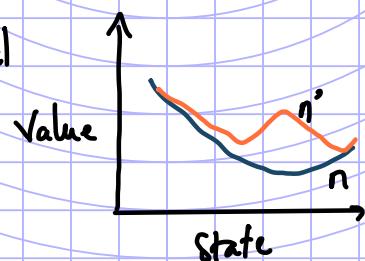
$$v_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_k(s')]$$

Control refers to the task of improving a policy

Policy improvement theorem

$$\pi'(s) = \operatorname{argmax}_a \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_\pi(s')]$$

$\pi' > \pi$ unless π is optimal



→ A synchronous methods
vs. asynchronous

1

4

8

16

$$\frac{1}{4} (x-1) + \frac{1}{4} (-2) = n$$