

- January 22, 2021 Notes: Chapter 5, Dynamic Programming

- Partially Observable Markov Decision Process (POMDP)

- ↳ These are more common in real world

- ↳ Two notions of state

- ① Environment can be viewed as a state

- ↳ We don't know the factors driving this. We view accumulated effects.

- ↳ $s_t^{(e)}$

- ② State you need to decide on next action.

- ↳ $s_t^{(a)}$

- ↳ Previously we assumed $s_t^{(e)} = s_t^{(a)} = s_t$, but this is not realistic

- ↳ More generally, You see observations o_t and infer state $s_t^{(e)}$ based on history of observations

- ↳ Observation function:

Modelled based on
understanding of business problem. Defined a priori.

$$z(s', a, o) = P[o_{t+1} = o \mid s_{t+1} = s', a_t = a]$$

- ↳ You still have probabilities P_A specified in POMDP

- ↳ MDP is special case of POMDP where:

$$o_t = s_t = s_t^{(e)} = s_t^{(a)}$$

- ↳ POMDP "guesses" s_t by maintaining Belief State

- ↳ See Slide 30 of Chapter 2 slides

- ↳ Belief states satisfy Markov property, but history must be recorded.
- ↳ Belief state of POMDP can be used as a single state of a giant MDP, which can be solved.
- ↳ Not practical, since state space is huge.
- ↳ Alternatively, treat observations as states.
- ↳ There are advanced techniques to solve POMDP, which lies beyond scope of class

* No assignments or exam questions on POMDP's; this is a current area of research.

- Chapter 2 Takeaways:

- MDP Bellman Policy Equations V^*
- MDP Bellman Optimality Equation V^*
- There exists an optimal policy, and each optimal policy satisfies optimal value function

- Chapter 3: Dynamic Programming

- Focus on proofs in textbook
- DP for both prediction and control problems
 - ↑ Predict V^*
 - ↑ Predict v^*
- Planning vs. Learning
- Bellman referred to DP as the Principle of Optimality
 - ↳ In CS, DP is used more broadly to mean "recursive algorithms with overlapping subproblems"
 - ↳ We refer to DP as "Algorithms for prediction and control."

- We will cover 3 DP algorithms based on the concept of Fixed Point

→ Fixed point of a function $f: X \rightarrow X$ is a point $x \in X$ such that $x = f(x)$

→ E.g. $f(x) = \cos(x)$, fixed point: $x^* = \cos(x^*)$

→ Applying a convergent function over and over again to solve this.

→ Banach Fixed-Point Theorem

Properties satisfied by physical distances
↓

X is non-empty set equipped with a complete metric $d: X \times X \rightarrow \mathbb{R}$.

Think of d as a distance
Let $f: X \rightarrow X$ be such that there exists an $L \in [0, 1)$ such that $d(f(x_1), f(x_2)) \leq L \cdot d(x_1, x_2)$ for all $x_1, x_2 \in X$

Then :

There exists a unique fixed point $x^* \in X$

* Proof of this theorem is beyond scope of course, but we will use this theorem

• Policy Evaluation (Prediction)

• MDP with $S = \{s_1, s_2, \dots, s_n\}$, $N = \{s_1, s_2, \dots, s_m\}$,

• Policy π

• Goal: find iterative algorithm to solve for π -implied MRP value function iteratively for large data sets

→ Evaluate Bellman Equation

↳ Define Bellman Policy operator that serves as
in Banach fixed-point theorem.

$$B^{\pi} : \mathbb{R}^m \rightarrow \mathbb{R}^m$$

↑
Vector representing a value function

$$\therefore B^{\pi}(v) = R^{\pi} + \gamma P^{\pi} \cdot v \quad \text{for any Value Function vector } v \in \mathbb{R}^m$$

↳ This can be re-expressed using v^{π}

$$\therefore v^{\pi} = B^{\pi}(v^{\pi})$$

↑
This is a fixed point!

- Metric $d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ defined as L^{∞} norm.

↳ Use this to prove B^{π} is a contraction function under L^{∞} norm

$$d(x, y) = \|x - y\|_{\infty} = \max_{s \in N} |(x - y)(s)|$$

↳ See slides 6-7 for proof. It takes role of L in Banach fixed point theorem to prove B^{π} is a contraction

- Policy Evaluation Algorithm: Start with any Value Function $v_0 \in \mathbb{R}^m$

↳ Iterate over $i = 0, 1, 2, \dots$, calculate in each iteration:

$$v_{i+1} = B^{\pi}(v_i) = R^{\pi} + \gamma P^{\pi} \cdot v_i$$

- ↳ Stop within tolerance
- ↳ Lower γ speeds up algorithm
- ↳ Runtime of each iteration is $O(m^2)$
- ↳ Constructing MRP from MDP or policy is $O(n^2 k)$, where k is the number of actions

- Greedy Policy: MDP Control Problem

↳ "Greedy" policy improvement technique

↳ Greedy Policy Function:

$$G : \mathbb{R}^m \longrightarrow \underbrace{(N \rightarrow A)}_{\text{Deterministic policy}}$$

$$\underbrace{G(V)(s)}_{\substack{\text{gives policy} \\ \text{gives action}}} = \pi^*(s) = \underbrace{\arg\max}_{a \in A} \left\{ R(s, a) + \gamma \sum_{s' \in N} P(s, a, s') \cdot V(s') \right\}$$

• Value Function Comparison: We say $X \geq Y$ iff $X(s) \geq Y(s)$ for all s

$$V^{\pi^*} = V^{G(V^*)} \geq V^*$$

We need to prove this.

• Policy Improvement Theorem:

• Policy \rightarrow Value Function \rightarrow Greedy Policy

\downarrow
Value function of
Greedy Policy

- Review slide 10, around 1 hr mark of lecture for prob. Prove By Induction

↳ End up with a tower of slightly improved value functions that converges on V^{π_0}

↳ Base case, note:

$$\begin{aligned} B^{\pi_0}(V^\pi)(s) &= \max_{a \in A} \left\{ R(s, a) + \gamma \sum_{s' \in S} p(s, a, s') V^\pi(s') \right\} \\ &\text{↳ greedy policy, deterministic} \\ &= \max_{a \in A} Q^\pi(s, a) \end{aligned}$$

↳ $V^\pi(s)$ is a weighted average over $Q^\pi(s, \cdot)$, while $B^{\pi_0}(V^\pi)(s)$ is a maximum value, which will always be greater than or equal to weighted avg.

↳ For induction step, via Monotonicity of B^π operator (for any π)

$$B^\pi : X \geq Y \Rightarrow B^\pi(X) \geq B^\pi(Y)$$

↳ Prove this!

↳ Use monotonicity problem to show
 $(B^{\pi_0})^{l+1}(V^\pi) \geq (B^{\pi_0})^l(V^\pi)$

* Review lecture from slide 12, 1 hr, 10 min mark

- Policy Iteration Algorithm: Start with Value Function $V_0 \in \mathbb{R}^n$

- Iterate over $j = 0, 1, 2, \dots$, calculate for each iteration:

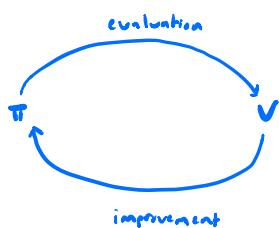
$$\text{Deterministic Policy} \quad \pi_{jH} = G(v_j)$$

$$\text{Value Function} \quad v_{jH} = \lim_{t \rightarrow \infty} (B^{\pi_{jH}})^t(v_j)$$

- Continue until v_j converges

* Iterate between Value Function improvement and Policy Evaluation

→ See slide 15.



→ Running Time : $O(n^2 k)$ ← How do you prove this?

- Value Function Iteration:

→ Bellman Optimality Operator

→ Tweak definition of Greedy Policy function ($\arg\max$ to \max)

→ Bellman Optimality Operator $B^* : \mathbb{R}^m \rightarrow \mathbb{R}^m$ defined as:

$$B^*(v)(s) = \max_{a \in A} \left\{ R(s, a) + \gamma \sum_{s' \in S} P(s, a, s') \cdot v(s') \right\}$$

• Review Slide 17

V^* is fixed point of B^*

$$\hookrightarrow V^* = B^*(V^*)$$

\hookrightarrow Prove that V^* is a contraction in next lecture

\hookrightarrow Beginning on Slide 19