# CS292F StatRL Lecture 2 Markov Decision Processes

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

# Recap: Markov Decision processes (MDP) parameterization

rolled out

$$\tau = (S_1, A_1, R_1, S_2, A_2, R_2, \dots)$$

- Infinite horizon / discounted setting

$$\mathcal{M}(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$$

$S_1 \sim \mu_0 \quad A_1 \sim \pi(a \mid S_t = S_1)$

$S_2 \sim P(s \mid S_t = S_1, A_t \in A_1)$

Transition kernel:
$$P : S \times A \to \Delta(S) \quad i.e. \quad P(S' \mid S, a)$$

(Expected) reward function:
$$r : S \times A \to \mathbb{R} / [0, R_{max}]$$

$[0, 1]$

$$\mathbb{E}[R_t \mid S_t = S, A_t = a] =: r(s, a)$$

Initial state distribution
$$\mu_0 \in \Delta(S)$$

Discounting factor: $0 \leq \gamma \leq 1$

Horizon $\dfrac{1}{1 - \gamma} = 1 + \gamma + \gamma^2 + \dots$

# Recap: Reward function and Value functions

- Immediate reward function r(s,a,s')
  - expected **immediate** reward

  $$r(s, a, s') = \mathbb{E}[R_1 | S_1 = s, A_1 = a, S_2 = s']$$

  $$r^\pi(s) = \mathbb{E}_{a \sim \pi(a|s)}[R_1 | S_1 = s]$$

- state value function: $V^\pi(s)$
  - expected **long-term** return when starting in *s* and following $\pi$

  $$V^\pi(s) = \mathbb{E}_\pi[R_1 + \gamma R_2 + ... + \gamma^{t-1} R_t + ... | S_1 = s]$$

- state-action value function: $Q^\pi(s,a)$
  - expected **long-term** return when starting in *s*, performing *a*, and following $\pi$

  $$Q^\pi(s, a) = \mathbb{E}_\pi[R_1 + \gamma R_2 + ... + \gamma^{t-1} R_t + ... | S_1 = s, A_1 = a]$$

# Recap: Optimal value function and the MDP planning problem

$$V^\star(s) := \sup_{\pi \in \Pi} V^\pi(s)$$

$$Q^\star(s,a) := \sup_{\pi \in \Pi} Q^\pi(s,a).$$

Goal of MDP planning:

$$\text{Find } \pi^* \text{ such that } V^{\pi^*}(s) = V^*(s) \quad \forall s$$

Approximate solution:

$$\pi \text{ is } \epsilon\text{-optimal if } V^\pi \geq V^*(s) - \epsilon\mathbf{1}$$

# Recap: General policy, Stationary policy, Deterministic policy

- General policy could depend on the entire history

$$\pi : (\mathcal{S} \times \mathcal{A} \times \mathbb{R})^* \times \mathcal{S} \to \Delta(\mathcal{A})$$

memoryless

- Stationary policy

$$\pi : \mathcal{S} \to \Delta(\mathcal{A})$$

- Stationary, Deterministic policy

$$\pi : \mathcal{S} \to \mathcal{A}$$

# Recap: We showed the following results about MDPs.

- **Proposition:** It suffices to consider stationary policies.
  1. Occupancy measure

$$V_\mu^\pi(s) = \sum_{t=1}^{\infty} \gamma^{t-1} \cdot d^\pi(S_t = s)$$

$$V^\pi(\theta) = \langle V_{d}^{\pi}(s,a), \gamma(s,a) \rangle$$

$$V_\mu^\pi(s,a) = \sum_{t=1}^{\infty} \gamma^{t-1} \cdot d^\pi(S_t = s, A_t = a)$$

$$\exists \; \pi' \; \text{stationary} \; s.t. \; V^\pi(s,a) = V^{\pi'}(s,a)$$

  2. There exists a stationary policy with the same occupancy measure

- **Corollary:** There is a stationary policy that is optimal for all initial states.
  - Proof sketch: 1. Construct an optimal non-stationary policy. 2. Apply the above proposition.

# Bellman equations – the fundamental equations of MDP and RL

- For stationary policies there is an alternative, recursive and more useful way of defining the V-function and Q function

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V^\pi(s')] = \sum_a \pi(a|s) Q^\pi(s,a)$$

*(handwritten annotations:)* $E_{a \sim \pi(a|s)}$  $E_{s' \sim P(\cdot|s,a)}$  immediate reward  discounted future reward under $\pi$  $Q^\pi(s,a)$

- Exercise:
  - Prove Bellman equation from the (first principle) definition.

  - Write down the Bellman equation using Q function alone.

$$Q^\pi(s,a) = ? \quad \sum_{s'} P(s'|s,a)\left[ r(s,a,s') + \gamma \sum_{a'} \pi(a'|s') Q^\pi(s',a') \right]$$

*(handwritten annotation:)* $V^\pi(s')$

# Deriving Bellman Equation for stationary policies

$$V^{\pi}(s) = \mathbb{E}^{\pi}\left[\sum_{t=1}^{\infty} \gamma^{t-1} r(S_t, A_t) \Big| S_1 = s\right]$$

Law of Total Expectation
$$= \mathbb{E}^{\pi}\left[r(S_1, A_1) \Big| S_1 = s\right] + \mathbb{E}^{\pi}\left[\sum_{t=2}^{\infty} \gamma^{t-1} r(S_t, A_t) \Big| S_2 = s'\right] \quad \overset{u}{t} = t-1$$

$$\underbrace{\sum_{S_2} P^{\pi}(S_2 = s' | S_1 = s)}$$

$$= r^{\pi}(s) + \gamma \sum_{S_2} P^{\pi}(S_2 = s' | S_1 = s) \, \mathbb{E}^{\pi}\left[\sum_{\overset{u}{t}=1}^{\infty} \gamma^{\overset{u}{t}-1} r(S_{\overset{u}{t}}, A_{\overset{u}{t}}) \Big| S_1 = s'\right]$$

By stationary
$$= r^{\pi}(s) + \gamma \sum_{S_2} P^{\pi}(S_2 = s' | S_1 = s) \, V^{\pi}(s')$$

$$\overset{R^S}{\underset{v}{\phantom{.}}} \qquad \overset{(P^{\pi})}{\phantom{.}}$$

$$\underline{V^{\pi} = r^{\pi} + \gamma \, P^{\pi} \cdot V^{\pi}}$$

$$\in R^{S \times S}$$

$$P^{\pi}(s' | s) = \sum_{a} P(s' | s, a) \cdot \pi(a | s)$$
$$\underset{R^{S \times S}}{\phantom{.}} \qquad \underset{R^{S \times SA}}{\phantom{.}}$$

$$P^{\pi} = (\text{Transition Matrix})^{\top}$$

Law of total expectation
$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]]$$

8

# Bellman equations in matrix forms

- Lemma 1.4 (Bellman consistency): For stationary policies, we have

$$V^\pi(s) = Q^\pi(s, \pi(s)). = \mathbb{E}\left[Q^\pi(s,a)\right]$$

$$a \sim \pi(a|s)$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}\left[V^\pi(s')\right].$$

- In matrix forms:

Marginalized $\in \mathbb{R}^{S \times S}$

$$V^\pi = r^\pi + \gamma P^\pi V^\pi \iff \left(I - \gamma P^\pi\right) V^\pi = r^\pi \in \mathbb{R}^S$$

$$Q^\pi = r + \gamma P V^\pi$$

point over $a'$

$$Q^\pi = r + \gamma \tilde{P}^\pi Q^\pi. \qquad \left(I - \gamma P^\pi\right) Q^\pi = r \in \mathbb{R}^{SA}$$

$\in \mathbb{R}^{SA \text{ by } SA}$

# Closed-form solution for solving for value functions

$$V^{\pi} = r^{\pi} + \gamma P^{\pi} V^{\pi}$$

$$Q^{\pi} = r + \gamma P V^{\pi}$$

$$Q^{\pi} = r + \gamma P^{\pi} Q^{\pi} \ .$$

$$V^{\pi} = \left( I - \gamma P^{\pi} \right)^{-1} \gamma \bar{r}$$

$$V^{\pi}(\mu) = \sum_{s,a} \gamma(s,a) \, \nu_{\mu}^{\pi}(s,a)$$

$$= \langle r, \nu_{\mu}^{\pi} \rangle$$

# Duality between *value functions* and *occupancy measures*

$$V^\pi = r^\pi + \gamma P^\pi V^\pi$$

$$Q^\pi = r + \gamma P V^\pi$$

$$Q^\pi = r + \gamma \breve{P}^\pi Q^\pi \,.$$

$$\begin{cases} V^\pi = \left(I - \gamma P^\pi\right)^{-1} \gamma^\pi \\[2mm] \nu^\pi = \left(I - \gamma (P^\pi)^\top\right)^{-1} \mu \end{cases}$$

$$\nu^\pi(s) = \mu(s) + \gamma \sum_{s'} \nu^\pi(s') \cdot p^\pi(s|s')$$

$$\nu^\pi = \mu + \gamma (P^\pi)^\top \nu^\pi$$

$$\nu^\pi(s,a) = \underbrace{\mu(s) \cdot \pi(a|s)}_{\mu^\pi(s,a)} + \gamma \sum_{s'} \nu^\pi(s') \left(\sum_{a'} p^\pi(s|s',a') \pi(a'|s')\right)$$

$$\pi(a|s)$$

$$\breve{\nu}^\pi = \mu^\pi + \gamma (\breve{P}^\pi)^\top \breve{\nu}^\pi$$

$$\nu^\pi(s,a) = \mu^\pi(s,a) + \gamma \sum_{s'} \sum_{a'} \nu^\pi(s',a') \underbrace{p^\pi(s,a|s',a')}_{R^{SA \times SA}}$$

# Invertibility of the matrix $I - \gamma P^\pi$

$(A \text{ is full rank} \iff A^\top \text{ is full rank})$

Corollary 1.5 in AJKS: the matrix $I - \gamma P^\pi$ is full rank / invertible for all gamma < 1.

Proof:

identity

$$\|(I - \gamma P^\pi)x\|_\infty = \|x - \gamma P^\pi x\|_\infty$$

triangular inequality $\to$

$$\geq \|x\|_\infty - \gamma \|P^\pi x\|_\infty$$
$$\geq \|x\|_\infty - \gamma \|x\|_\infty$$
$$= (1 - \gamma)\|x\|_\infty > 0$$

$\gamma < 1$

$\sum_s \leq 1 \quad \sum_s = 2$

$p(s'|s) \sim \sum_s \frac{1}{\leq} |S|$

$P^{\pi} \quad \leftarrow S$

$P^{\pi}x = \begin{pmatrix} \langle P^{\pi}[1,:], x \rangle \\ \langle P^{\pi}[2,:], x \rangle \\ \vdots \end{pmatrix}$

$\langle P^{\pi}[i,:], x \rangle$
$\leq \|P^{\pi}[i,:]\|_1 \|x\|_\infty$
$\leq \|x\|_\infty$

12

# Bellman optimality equations characterizes the optimal policy

expected immediate reward

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

discounted future reward

by the Optimal policy

- system of n non-linear equations
- solve for V*(s)
- easy to extract the optimal policy

- having Q*(s,a) makes it even simpler

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

# Proposition: There is a *deterministic, stationary* and *optimal* policy.

- And it is given by:

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

- Proof:

$\pi^*$ is stationary

$$V^{\pi'}(s) \leq V^*(s) = V^{\pi^*}(s) = \mathbb{E}_{a \sim \pi^*(a|s)} Q^{\pi^*}(s,a) \leq \max_a Q^{\pi^*}(s,a)$$

$\forall \pi$

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

Substitute $\pi = \pi'$

define $\pi'(s) = \arg\max_a Q^*(s,a)$

$\pi'$ is stationary
$\pi'$ is deterministic

$V^{\pi'}(s)$

# The crux of solving the MDP planning problem is to construct Q*

- In the remainder of this lecture, we will talk about two approaches

  1. By solving a Linear Program

  2. By solving Bellman equations / Bellman optimality equations.

# The linear programming approach

- Solve for V* by solving the following LP

$$\min_{V \in \mathbb{R}^S} \quad \sum_s \mu(s) V(s)$$

$$\text{subject to} \quad V(s) \geq r(s,a) + \gamma \sum_{s'} P(s'|s,a) V(s') \quad \forall a \in \mathcal{A},\ s \in \mathcal{S}$$

$V^{\pi^*}(\mu)$

$(Ye, 1990s)$

Substitue $U = U^*$ $\quad \sum_s \mu(s) U^*(s) = V^*(\mu)$

$-) \quad V(s) \geq \max_a r(s,a) + \gamma \sum_{s'} P(s'|s,a) V(s')$

# The linear programming approach

- Solve for V* by solving the following LP

$$\min \quad \sum_s \mu(s)V(s)$$

$$\text{subject to} \quad V(s) \geq r(s,a) + \gamma \sum_{s'} P(s'|s,a)V(s') \quad \forall a \in \mathcal{A}, s \in \mathcal{S}$$

**Quiz 1:** Once we have V*, how to construct Q*?

$$\pi^*(s) = \underset{a}{\arg\max} \, Q^*(s,a)$$

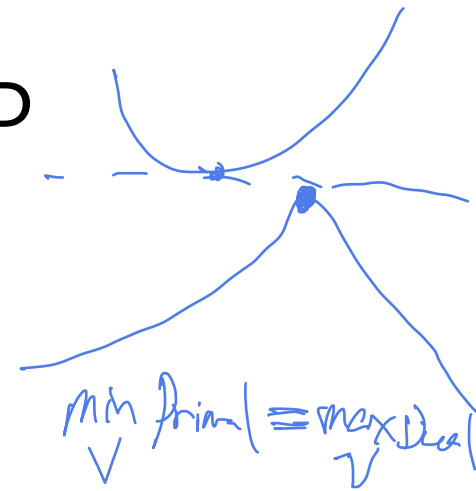$$Q^*(s,a) = [r(s,a) + \gamma \sum P(s'|s,a) \cdot V^*(s')]$$

16

# The Lagrange dual of the LP

$$\max_{\nu} \quad \sum_{s,a} \nu(s,a) r(s,a)$$

subject to

$$\nu \geq 0$$

$$\sum_{z} \nu(s,a) = \mu(s) + \gamma \sum_{s',a'} P(s|s',a') \nu(s',a')$$

- Exercise: Deriving the dual by applying the standard procedure.

# The Lagrange dual of the LP

$$\max_{\nu} \quad \sum_{s,a} \nu(s,a) r(s,a)$$

$$\text{subject to} \quad \nu \geq 0$$

$$\sum_{z} \nu(s,a) = \mu(s) + \gamma \sum_{s',a'} P(s|s',a') \nu(s',a')$$

$\nu \in \mathbb{R}^{SA}$

- Exercise: Deriving the dual by applying the standard procedure.

**Quiz 2:** Once we have the solution how to construct the policy?

$\pi^*(a|s) = \dfrac{\nu^{\pi^*}(s,a)}{\sum_a \nu^{\pi^*}(s,a)}$

$\nu^*(s,a) = \nu^{\pi^*}(s,a) = \nu^{\pi^*}(s) \cdot \pi^*(a|s)$

# Value iterations for MDP planning

- Recall: Bellman optimality equations

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

$$Q(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \max_{a' \in \mathcal{A}} Q(s',a') \right].$$

$$\mathcal{T}Q = r + PV_Q \qquad \text{where} \qquad V_Q(s) := \max_{a \in \mathcal{A}} Q(s,a).$$

**Theorem** 1.8 (AJKS):  Q = Q* if and only if Q satisfies the Bellman optimality equations.

# Value iterations for MDP planning

- The value iteration algorithm iteratively applies the Bellman operator until it converges.

  1. Initialize $Q_0$ arbitrarily      $Q_0 = 0$

  2. for i in 1,2,3,…, k,  update $\quad Q_i = \mathcal{T}Q_{i-1}$

  3. Return $Q_k$

# Value iterations for MDP planning

- The value iteration algorithm iteratively applies the Bellman operator until it converges.

  1. Initialize $Q_0$ arbitrarily

  2. for i in 1,2,3,…, k, update $Q_i = \mathcal{T} Q_{i-1}$

  3. Return $Q_k$

- **What is the right question to ask here?**

1. $\lim_{k \to \infty} Q_k = Q^*$ ?

2. $\| Q_k - Q^* \|_\infty \leq \epsilon(k)$   rate of convergence

3. Iterative Complexity: $\epsilon$ as an input $k \geq func(\epsilon)$

# Convergence analysis of VI

$$\mathcal{T}Q = r + \gamma P \cdot V_Q, \quad V_Q = \max_a Q(\cdot, a)$$

- Lemma 1. The Bellman operator is a γ-contraction.

*For any two vectors $Q, Q' \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$,*

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \boxed{\gamma \|Q - Q'\|_\infty}$$

$\sum_j = 1$

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty = \gamma \|PV_Q - PV_{Q'}\|_\infty = \gamma \|P(V_Q - V_{Q'})\|_\infty$$

operator
norm of
$P$ in $\ell\infty$

$\longrightarrow \leq \gamma \|V_Q - V_{Q'}\|_\infty = \gamma \max_s |V_Q(s) - V_{Q'}(s)| \leq \gamma \max_{s,a} |Q(s,a) - Q'(s,a)|$

① if $V_Q(s) \geq V_{Q'}(s)$ for those $s$

$\gamma \max_s (V_Q(s) - V_{Q'}(s)) \leq \gamma \, Q(s,a) - \max_a Q'(s,a)$

$a = \arg\max_a Q(s,a)$

$\leq \gamma \, Q(s,a) - Q'(s,a)$

$\leq \gamma |Q(s,a) - Q'(s,a)|$

20

# Convergence analysis of VI

$$\|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

- Lemma 2. Convergence of the Q function.

$$Q' = Q^*, \qquad \boxed{TQ^* = Q^*}$$

$$\|Q_k - Q^*\|_\infty = \|T Q_{k-1} - Q^*\|_\infty \leq \gamma \|Q_{k-1} - Q^*\|_\infty$$

$$\|Q_0 - Q^*\|_\infty \leq \frac{1}{1-\gamma}$$

$$\cdots \leq \gamma^k \cdot \frac{1}{1-\gamma} = \frac{(1-(1-\gamma))^k}{1-\gamma} < \frac{e^{-(1-\gamma)k}}{1-\gamma}$$

$$0 \leq r(s,a) \leq 1 \qquad \left| \sum_{t=1}^{\infty} \gamma^{t-1} r(s,a) \right| \leq \frac{1}{1-\gamma}$$

**Quiz 3**: Computing "Iteration complexity" from "convergence bound"?

$$\varepsilon = \frac{e^{-(1-\gamma)k}}{1-\gamma} \iff k = \frac{\log \frac{1}{\varepsilon(1-\gamma)}}{1-\gamma}$$

$$\lim_{n \to \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1}$$

$$\left(1 - \frac{1}{n}\right)^n \leq e^{-1} \text{ for all } n \geq 1$$

# Convergence of the Q function implies the convergence of the value of the *induced* policy.

$\pi_Q(s) = \arg\max_a Q(s,a)$

---

Lemma 1.11 AJKS (Q-error amplification):

$$V^{\pi_Q} \geq V^\star - \frac{2\|Q - Q^\star\|_\infty}{1 - \gamma} \mathbb{1}.$$

---

**Proof:** Fix state $s$ and let $a = \pi_Q(s)$. We have:

$$
\begin{aligned}
V^\star(s) - V^{\pi_Q}(s) &= Q^\star(s, \pi^\star(s)) - Q^{\pi_Q}(s, a) \\
&= Q^\star(s, \pi^\star(s)) - Q^\star(s, a) + Q^\star(s, a) - Q^{\pi_Q}(s, a) \\
&= Q^\star(s, \pi^\star(s)) - Q^\star(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}[V^\star(s') - V^{\pi_Q}(s')] \\
&\leq Q^\star(s, \pi^\star(s)) - Q(s, \pi^\star(s)) + Q(s, a) - Q^\star(s, a) \\
&\quad + \gamma \mathbb{E}_{s' \sim P(s,a)}[V^\star(s') - V^{\pi_Q}(s')] \\
&\leq 2\|Q - Q^\star\|_\infty + \gamma \|V^\star - V^{\pi_Q}\|_\infty.
\end{aligned}
$$

where the first inequality uses $Q(s, \pi^\star(s)) \leq Q(s, \pi_Q(s)) = Q(s, a)$ due to the definition of $\pi_Q$.

# An alternative method: policy iteration

Initialize a policy $\pi_0$ arbitrarily.
for k= 1,2,3,4,...

1. *Policy evaluation.* Compute $Q^{\pi_k}$

2. *Policy improvement.* Update the policy: $\pi_{k+1} = \pi_{Q^{\pi_k}}$

*Solution to Bellman equation for $\pi$*

$$Q^{\pi} = (I - \gamma P^{\pi})^{-1} \cdot r$$

**Theorem 1.14.** *(Policy iteration convergence). Let $\pi_0$ be any initial policy. For $k \geq \frac{\log \frac{1}{(1-\gamma)\epsilon}}{1-\gamma}$, the k-th policy in policy iteration has the following performance bound:*

$$Q^{\pi^{(k)}} \geq Q^{\star} - \epsilon \mathbb{1} .$$

# Computational complexity of these MDP solvers

- VI:

$$S^2 \cdot A \cdot \frac{\log \frac{1}{(1-\gamma)^2 \varepsilon}}{1-\gamma}$$

$$\underset{\text{apply } \tau}{\uparrow}$$

$$\varepsilon = 0$$

- PI:

$$|SA|^3 \cdot \frac{\log \frac{1}{(1-\gamma)\varepsilon}}{1-\gamma} \implies (S^3 + S^2 A) \cdot \frac{\log \frac{1}{(1-\gamma)\varepsilon}}{1-\gamma}$$

- LP:

$$\text{poly } (S, A)$$

# Strongly polynomial algorithms are independent to ε

*a version of SIMPLEX method*

| | Value Iteration | Policy Iteration | LP-Algorithms |
|---|---|---|---|
| Poly? | $\|\mathcal{S}\|^2\|\mathcal{A}\|\frac{L(P,r,\gamma)\log\frac{1}{1-\gamma}}{1-\gamma}$ | $(\|\mathcal{S}\|^3+\|\mathcal{S}\|^2\|\mathcal{A}\|)\frac{L(P,r,\gamma)\log\frac{1}{1-\gamma}}{1-\gamma}$ | $\|\mathcal{S}\|^3\|\mathcal{A}\|L(P,r,\gamma)$ |
| Strongly Poly? | ✗ | $(\|\mathcal{S}\|^3+\|\mathcal{S}\|^2\|\mathcal{A}\|)\cdot\min\left\{\frac{\|\mathcal{A}\|^{\|\mathcal{S}\|}}{\|\mathcal{S}\|},\frac{\|\mathcal{S}\|^2\|\mathcal{A}\|\log\frac{\|\mathcal{S}\|^2}{1-\gamma}}{1-\gamma}\right\}$ | $\|\mathcal{S}\|^4\|\mathcal{A}\|^4\log\frac{\|\mathcal{S}\|}{1-\gamma}$ |

$(\|\mathcal{S}\|)^2$

# Next lecture

- Approximate / randomized solvers for MDP

- MDP / RL with generative models