# Unit 1: MDP basics

## NDP formulations



Agent

action            Feedback

Enviroment

## Markov Decision process (MDP)

focus on this
for simplicity

Temporal correlation:     actions has influence on the futures

$\Big\{$

Infinite-horizon (discounted) MDP:     $M = (S, A, P, r, \gamma)$

Finite-horizon MDP:     $M = (S, A, \{P_h\}_{h \in [H]}, \{r_h\}_{h \in [H]}, H)$

①

## Interaction protocol: (Discounted MDP)

- Start: $s_0 \sim M$
- for time step $t = 1, 2, \ldots, T$:
  - agent $\Big\{$ takes $a_t \in A$

    obtains rewards: $r_t = r(s_t, a_t)$

    observes next state: $s_{t+1} \sim P(\cdot \mid s_t, a_t)$

$\Big\{$

$S$: State space

$A$: action space

$r: S \times A \longrightarrow [0,1]$ reward function

$P: S \times A \longrightarrow \Delta(S)$, transition probability

$\gamma$: discount factor          $P(s' \mid s, a)$

②

## Enviroment:

$$P, r$$

$$\begin{cases} ① \quad \text{known} \quad \rightarrow \text{Dynamic Programing} \\\\ ② \quad \text{Simulator} \quad \rightarrow \text{"generative model"} \\\\ ③ \quad \text{unknown, have to play from beginning} \\ \qquad \qquad \uparrow RL \end{cases}$$

## Policy:

History - dependent policies: $\Pi^{hist} = \left\{ (S \times A \times \mathbb{R})^{*} \times S \longrightarrow \Delta(A) \right\}$

Stationary policies: $\Pi^{stn} = \left\{ S \longrightarrow \Delta(A) \right\}$

Deterministic policies: $\Pi^{det} = \left\{ S \longrightarrow A \right\}$

③

## Value:

$V_M^\pi(s)$ : cumulative rewards

$$V_M^\pi(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t\, r(s_t, a_t) \,\Big|\, \pi,\, s_0 = s\right]$$

$$Q_M^\pi(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t\, r(s_t, a_t) \,\Big|\, \pi,\, s_0 = s,\, a_0 = a\right]$$

## Goal:

Find an optimal policy $\pi^*$, i.e.,

$$V_M^{\pi^*}(s) \geqslant V_M^\pi(s) \qquad \forall\, (s, \pi) \in (S, \pi^{hist})$$

④

Theorem:

$$\boxed{\exists\, \pi^* \in \Pi^{det}}$$

Bellman equations: $\forall \pi \in \Pi^{stn}$

$$
\begin{cases}
V_M^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ Q_M^\pi(s,a) \right] \\[3mm]
Q_M^\pi(s,a) = r(s,a) + \gamma \mathbb{E}_{\substack{a \sim \pi(\cdot|s) \\ s' \sim P(\cdot|s,a)}} \left[ V^\pi(s') \right]
\end{cases}
$$

Bellman optimality equations:

$$
\begin{cases}
V_M^{\pi^*}(s) = \max_{a \in A} Q_M^{\pi^*}(s,a) \\[3mm]
Q_M^{\pi^*}(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \max_{a'} Q_M^{\pi^*}(s',a') \right]
\end{cases}
$$

Bellman optimality operator: $\quad T: \{S \times A \to \mathbb{R}\} \longrightarrow \{S \times A \to \mathbb{R}\}$

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

5

# Unit 2:   MDP Planning

$\left( \text{consider} \quad \pi \in \Pi^{stn} \right)$

$\nearrow$ known $r, P$

Notations: .     If context is clear:

$$V_M^\pi \longrightarrow V^\pi, \quad Q_M^\pi \longrightarrow Q^\pi, \quad V_M^{\pi^*} \longrightarrow V^*, \quad Q_M^{\pi^*} \longrightarrow Q^*$$

. Vectorization everything.

$$V^\pi = \left[ V^\pi(s) \right]_{s \in S} \in \mathbb{R}^{|S|}, \quad Q^\pi \in \mathbb{R}^{|S| \cdot |A|}$$

$$r \in \mathbb{R}^{|S| \cdot |A|}, \quad P = \left[ P(s'|s,a) \right]_{s', s, a} \in \mathbb{R}^{|S| \cdot |A| \times |S|}$$

$$P^\pi = \left[ P^\pi_{(s',a'),(s,a)} \right] \in \mathbb{R}^{|S| |A| \times |S| \cdot |A|}$$

$$\nearrow P^\pi_{(s',a'),(s,a)} = P(s'|s,a) \, \pi(a'|s').$$

Greedy policy:     Given $Q \in \{ S \times A \to \mathbb{R} \}$

$$\pi_Q(s) \in \underset{a \in A}{\text{argmax}} \; Q(s,a), \; \forall s$$

# vectorized Bellman equations

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

claim :

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

proof : - Only need to prove $I - \gamma P^\pi$ is invertible

- $\forall x \in \mathbb{R}^{|S|}$, s.t. $\|x\|_\infty > 0$,

$$\| (I - \gamma P^\pi) x \|_\infty = \| x - \gamma P^\pi x \|_\infty$$
$$\geqslant \|x\|_\infty - \gamma \| P^\pi x \|_\infty$$
$$\geqslant \|x\|_\infty - \gamma \| x \|_\infty \qquad \because \text{row} (P^\pi x)$$
$$= (1 - \gamma) \|x\|_\infty > 0 \qquad \text{is an average of the coordinates of } x.$$

Claim: $(1-\gamma)\left[I - \gamma P^\pi\right]^{-1} = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t (P^\pi)^t$

<span style="color:green">$\underbrace{\qquad\qquad\qquad}$</span>

<span style="color:green">a mixture of distributions</span>

Claim (contraction):

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

$$\forall Q, Q' \in \{S \times A \to \mathbb{R}\}$$

## Value iteration:

- set $Q_0 = 0^{|S| \cdot |A|}$
- iterate: $Q_{t+1} = TQ_t$, $t = 0, 1 \ldots$

Claim

$$\| Q_t - Q^* \|_\infty \leq \gamma^t \| Q^* \|_\infty$$

## Policy iteration:

- Start from an arbitrary policy $\pi_0$

- For $t = 0, 1, \ldots,$ :
  - Policy evaluation: $Q^{\pi_t}$
  - Greedify: $\pi_{t+1} = {}^\pi Q^{\pi_t}$

Claim:
- $Q^{\pi_t} \leq TQ^{\pi_t} \leq Q^{\pi_{t+1}}$
- $\| Q^{\pi_{t+1}} - Q^* \|_\infty \leq \gamma \| Q^{\pi_t} - Q^* \|_\infty$

# Linear programming:

## Primal:

subject to
$$\left\{ \begin{array}{l} \min\limits_{V \in \mathbb{R}^{|S|}} V^T \mu \\ V(s) \geq r(s,a) + \gamma \sum\limits_{s'} P(s'|s,a) V(s') \quad \forall s \in S, a \in A \end{array} \right.$$

Claim: $V^*$ is the unique solution

## Dual LP:

$$\arg\max \frac{1}{1-\gamma} d^T r$$

subject to: $\sum\limits_{a \in A} d(s,a) = (1-\gamma) \mu(s) + \gamma \sum\limits_{s',a'} P(s|s',a') d(s',a')$

$$\forall s \in S$$

# Unit 3:    Simulator setting
## ( Generative models)

Simulator :



assume $r$ is known to learner

Goal:

. Find $\epsilon$-optimal value function $\widehat{Q}$ :     (value guarantee / bound)

$$\| Q^* - \widehat{Q} \|_\infty \leq \epsilon$$

. Find $\epsilon$-optimal policy $\widehat{\pi}$ :     (policy guarantee / bound)

$$\| Q^* - Q^{\widehat{\pi}} \|_\infty \leq \epsilon$$     ← more subtle

## Model-based method

- maximum likelihood estimate (MLE) of the transition kernel and use it as a "plug-in"

- For each $(s,a) \in S \times A$, draw $N$ samples of $s' \sim P(\cdot|s,a)$

- Let $\widehat{P}$ be the empirical model, i.e.

$$\widehat{P}(s'|s,a) = \frac{count \, (s', s, a)}{N}$$

- \# samples $= |S| \cdot |A| \cdot N$

## Notations

- Denote $\hat{P} = \left[\hat{P}(s'|s,a)\right]_{(s,a,s')} \in \mathbb{R}^{|S|\cdot|A| \times |S|}$

- Let $\hat{M} = (S, A, r, \hat{P}, \gamma, \mu)$

- Let $\hat{Q}^{\pi} = Q^{\pi}_{\hat{M}}$, $\quad \hat{\pi}$ is an optimal policy of $\hat{M}$

- Let $H = \dfrac{1}{1-\gamma}$ : (effective) horizon

# Coarse analysis

$\forall \varepsilon, \delta > 0$

**Theorem** if #samples $= |S|.|A|.N \gtrsim \dfrac{\gamma^2 H^4 |S|^2 |A| \log\left(\frac{|S|.|A|}{\delta}\right)}{\varepsilon^2}$

- w.p. a.l. $1-\delta$, $\|Q^* - \hat{Q}^*\|_\infty \leq \varepsilon$

- wpal $1-\delta$, $\|Q^* - Q^{\hat{\pi}}\|_\infty \leq \varepsilon$

## Proof:

$$\|Q^* - \hat{Q}^*\|_\infty = \|\max_\pi Q^\pi - \max_\pi \hat{Q}^\pi\|_\infty$$

$$\leq \max_\pi \|Q^\pi - \hat{Q}^\pi\|_\infty$$

uniform convergence

simulation lemma + concentration

$$\|Q^k - Q^{\hat{\pi}}\|_\infty = \|Q^* - \hat{Q}^* + \hat{Q}^* - Q^{\hat{\pi}}\|_\infty$$

$$\leq \|Q^* - \hat{Q}^*\|_\infty + \|\hat{Q}^{\hat{\pi}} - Q^{\hat{\pi}}\|_\infty$$

$$\leq 2 \cdot \max_\pi \|Q^\pi - \hat{Q}^\pi\|_\infty$$

## Simulation lemma:

$$Q^\pi - \hat{Q}^\pi = \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi$$

proof:     Recall:     $Q^\pi = r + \gamma P^\pi Q^\pi \Rightarrow \begin{cases} Q^\pi = (I - \gamma P^\pi)^{-1} r \\ r = (I - \gamma P^\pi) Q^\pi \end{cases}$

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

$$= (I - \gamma \hat{P}^\pi)^{-1} \left( (I - \gamma \hat{P}^\pi) Q^\pi - (I - \gamma P^\pi) Q^\pi \right)$$

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

$$= \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi$$

$$\| Q^\pi - \hat{Q}^\pi \|_\infty = \| \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi \|_\infty$$

$$\leq \frac{\gamma}{1-\gamma} \| (P - \hat{P}) V^\pi \|_\infty$$

$$= \frac{\gamma}{1-\gamma} \max_{(s,a)} \left| \left( P(\cdot|s,a) - \hat{P}(\cdot|s,a) \right)^T V^\pi(\cdot) \right|$$

Hölder
$$\leq \frac{\gamma}{1-\gamma} \max_{(s,a)} \| P(\cdot|s,a) - \hat{P}(\cdot|s,a) \|_1 \cdot \| V^\pi \|_\infty$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \max_{(s,a)} \| P(\cdot|s,a) - \hat{P}(\cdot|s,a) \|_1$$

Using McDiarmid's Inequality :

$$\| P(\cdot|s,a) - \hat{P}(\cdot|s,a) \|_2 \leq \sqrt{\frac{2 \log(e|S|.|A|/\delta)}{N}}$$

Holder's inequality :

$$\| P(\cdot | s,a) - \hat{P}(\cdot | s,a) \|_1 \leq \sqrt{|S|} \cdot \| P(\cdot | s,a) - \hat{P}(\cdot | s,a) \|_2$$

# Crude value bounds

- \# samples in the coarse analysis is linear in model complexity
- it comes from guarantees uniformly over all policies
- yet we only need to care about $\pi^*$ and $\hat{\pi}$.

- can \# samples be sublinear in model complexity?