

1) Infinite Horizon Discounted MDP

$$M = \{S, A, P, r, \gamma\}$$

S : space of possible states $s \in S$

A : space of possible actions $a \in A$

P : transition function $P: S \times A \rightarrow \Delta(S)$

r : reward function $r: S \times A \rightarrow \Delta(\mathbb{R})$

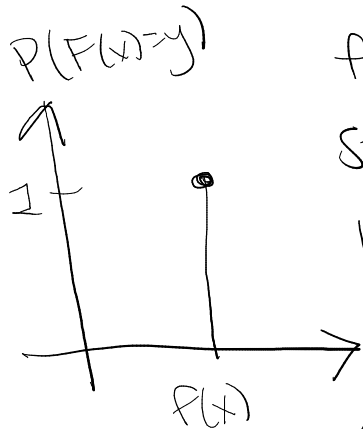
γ : discount factor $0 < \gamma < 1$

probability distributions

Aside: We can encode a deterministic

function $f: X \rightarrow Y$ as a stochastic one $F: X \rightarrow \Delta(Y)$

by $F(x) = f(x)$ w.p. 1.



sometimes as shorthand we will overload notation and write

e.g. $a = \pi(s)$ instead of $a \sim \pi(s)$ if the policy is deterministic.

Additionally, we will adopt the notation

$$F(y|x) = \mathbb{P}\{F(x) = y\}$$

(e.g. $\pi(a|s)$, $P(s'|a, s)$)

In this notation we can write the goal:

finding a policy $\pi: \mathcal{S} \rightarrow \mathcal{A}(\mathcal{S})$
that maximizes the (discounted)
cumulative reward.

$$\begin{array}{l} \text{maximize} \\ \pi \end{array} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \right. \\ \left. \begin{array}{l} s_{t+1} \sim P(s_t, a_t), s_0 \text{ given,} \\ a_t \sim \pi(s_t) \end{array} \right]$$

We will spend the semester learning how to solve this problem. In RL, we

do not assume that $P(\cdot, \cdot)$ is known, and therefore we have to solve the optimization using data.

For now, we suppose that P is known

2) Value and Q Function

allow us to reason about policy's long term effect.

$$V^{\pi}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, s_{t+1} \sim P(s_t, a_t), a_t \sim \pi(s_t) \right]$$

$$Q^{\pi}(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, s_{t+1} \sim P(s_t, a_t) \right. \\ \left. a_0 = a, a_t \sim \pi(s_t) \right]$$

Bellman Equations

Notice that

$$\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) = r(s_0, a_0) + \gamma \sum_{t=1}^{\infty} r(a_t, s_t)$$

$$(\text{let } t' = t+1) = r(s_0, a_0) + \gamma \sum_{t'=0}^{\infty} r(a_{t'+1}, s_{t'+1})$$

let's consider deterministic policies and reward functions.

This observation allows us to write

$$V^{\pi}(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} [V^{\pi}(s')]$$

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

Policy EV: what property of

expectations do we use?

How would the expressions change for stochastic reward functions and policies?

3) Policy Evaluation

How do we characterize how good a policy is? In terms of value function

Given MDP $\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, P, \gamma, r \}$ and policy π , what is V^π ?

function from $\mathcal{S} \rightarrow \mathbb{R}$

The Bellman equation:

$$\forall s, \quad V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} [V^\pi(s')]]$$

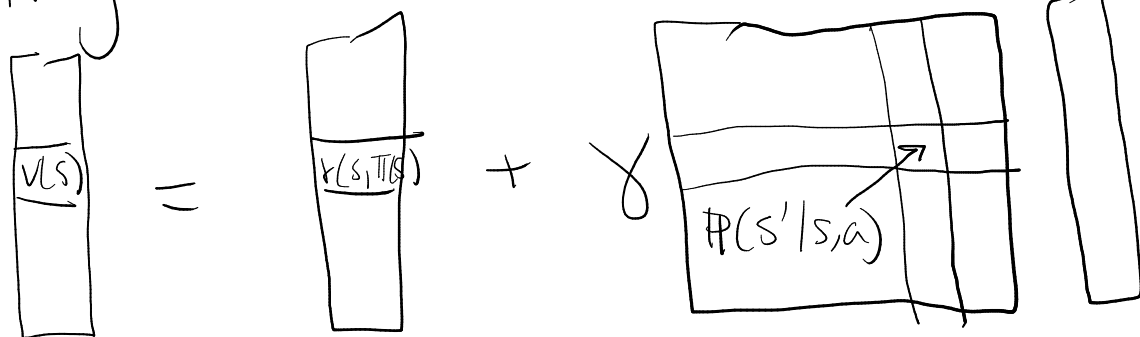
$$\forall s, \quad V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$$

denote $S = |\mathcal{S}|$

number of states

S linear constraints
 S unknowns

Writing in vector-matrix notation



$$V \in \mathbb{R}^S$$

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

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

$$V$$

Solving the linear equations

$$V = R + \gamma P V \longrightarrow V = (I - \gamma P)^{-1} R$$

This is valid as long as $I - \gamma P$ is invertible (HW0)

Exact solution! But $O(S^3)$ for matrix inversion...

4) Approximate Policy Evaluation

can we trade accuracy for faster computation?
Yes! Iterative Algorithm for fixed point.

Algorithm (Iterative PE)

initialize V^0

for $t=0, \dots$

$$V^{t+1} \leftarrow R + \gamma P V^t$$

Q: complexity per iteration?

A: matrix-vector multiply is $O(S^2)$

To show that this algorithm works, we will show a contraction, which is a general strategy for fixed point algorithms.

$$\text{Lemma: } \|V^{t+1} - V^\pi\|_\infty \leq \gamma \|V^t - V^\pi\|_\infty$$

Proof:

$$\begin{aligned} \|V^{t+1} - V^\pi\|_\infty &= \|R + \gamma P V^t - V^\pi\|_\infty \quad (\text{alg}) \\ &= \|R + \gamma P V^t - (R + \gamma P V^\pi)\|_\infty \quad (\text{Bellman}) \\ &= \gamma \|P(V^t - V^\pi)\|_\infty \end{aligned}$$

recall that each entry of this vector represents the expectation at index s , $|\mathbb{E}[V^t(s') - V^\pi(s')]|$
 $s' \sim P(s, \pi(s))$ (Jensen's)

$$\leq \mathbb{E} |V^t(s') - V^\pi(s')|$$

$$\text{so } \|P(V^t - V^\pi)\|_\infty \leq \max_s \mathbb{E} [|V^t(s') - V^\pi(s')|]$$

(expectation upper bounded by max)

$$\leq \max_{s'} |V^t(s') - V^\pi(s')| = \|V^t - V^\pi\|_\infty$$

$$\text{thus } \|V^{t+1} - V^\pi\|_\infty \leq \gamma \|V^t - V^\pi\|_\infty \quad \square$$

Theorem: after t iterations,

$$\|V^t - V^\pi\|_\infty \leq \gamma^t \|V^0 - V^\pi\|_\infty$$

i.e.,

$$\forall s, |V^t(s) - V^\pi(s)| \leq \gamma^t \max_s |V^0(s) - V^\pi(s)|$$

follows by repeated application of Lemma.

How many iterations necessary for ϵ accurate solution?

$$\gamma^t \|V^0 - V^\pi\|_\infty \leq \epsilon$$

$$\Rightarrow t \geq \log\left(\frac{\|V^0 - V^\pi\|_\infty}{\epsilon}\right) / \log(1/\gamma)$$

overall complexity

$$O(s^2 \log(1/\epsilon))$$

compare with $O(s^3)$ for exact.

5) State-Action Distribution

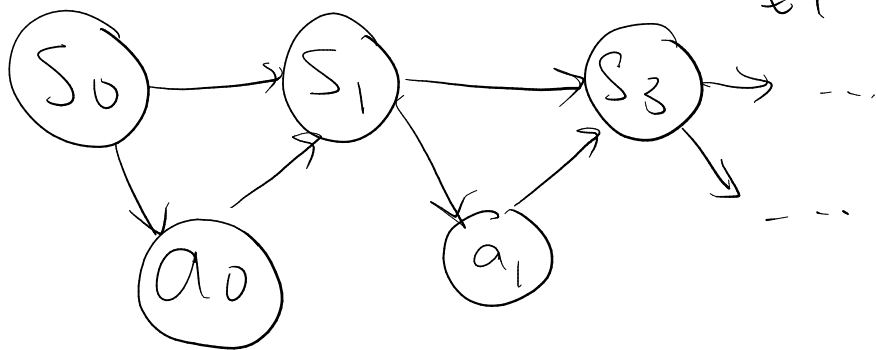
Trajectory of MDP up to step t :

$$(s_0, a_0, s_1, a_1, \dots, s_t, a_t)$$

What is the probability of a particular trajectory under policy π ?

considering possibly stochastic policies,

$$\begin{aligned} \mathbb{P}^\pi (s_0, a_0, \dots, s_t, a_t) &= \pi(a_0 | s_0) P(s_1 | s_0, a_0) \times \\ &\quad \pi(a_1 | s_1) P(s_2 | s_1, a_1) \times \dots \\ &\quad \times P(s_t | s_{t-1}, a_{t-1}) \pi(a_t | s_t) \end{aligned}$$



What is the probability of seeing (s, a) at timestep t , starting from s_0 ?

$$\mathbb{P}_t^\pi (s, a; s_0) = \sum_{\substack{a_{0:t-1}, \\ s_{0:t-1}}} \mathbb{P}^\pi (s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t = s, a_t = a)$$

Discounted Average State-Action Distribution

$$d_{s_0}^{\pi}(s, a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_n^{\pi}(s, a; s_0)$$

HW0: is this a valid distribution?

$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{s, a} d_{s_0}^{\pi}(s, a) r(s, a)?$$