

# Markov Decision Processes

- An MDP is a general way to model an online decision-making problem where any uncertain parameter is modelled in a **Bayesian manner** (i.e., as being drawn via some known stochastic process)
- MDPs can be defined over continuous spaces, and with continuous-time updates. We will focus (for now...) on **discrete time updates, and discrete (finite/countable) states** (This is sometimes called a **tabular MDP**)

- **Def:** A **Markov Chain** is a stochastic process  $(X_t)_{t=0}^{\infty}$  given by a **stochastic update** (i.e., randomized fn)

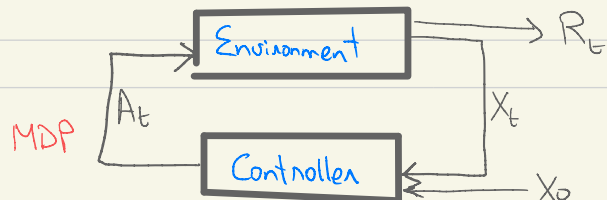
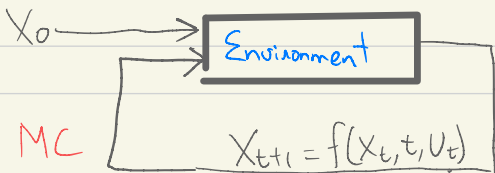
$$X_{t+1} = f(\overset{\text{state}}{X_t}, t, \underset{\text{independent n.o. ('disturbance')}}{U_t})$$

where  $U_1, U_2, \dots$  are iid  $U[0,1]$  r.v. (recall: ANY r.v.  $Y$  with cdf  $F$  can be constructed as  $Y = F^{-1}(U)$ )

- Any Markov chain comprises of the following 'inputs'

- **State space**  $S$
- **Initial state**  $X_0$  (or initial distr<sup>n</sup>  $\Pi_0$  over  $S$ )
- **(time  $t$ ) transition 'kernel'**  $P_t(x|X_t) = P[X_{t+1}=x|X_t]$

- An MDP interlaces a Markov chain with a 'control' module



Defn. A **Markov Decision Process** comprises of 3 interlacing sequences -

**States**  $X_0, X_1, X_2 \dots \in S$  (State space)  
**Actions**  $A_0, X_1, X_2 \dots \in A$  (Action space)  
**Rewards**  $R_1, R_2, \dots$  (Reward)

These are related via two functions

$$\begin{aligned}
 X_{t+1} &= f(X_t, A_t, U_t) \quad (\text{Transition function}) \\
 R_{t+1} &= g(X_t, A_t, U_t) \quad (\text{Reward function})
 \end{aligned}$$

## Notes

- The transitions can be represented via a **transition kernel**

$$T_t(x | X_t, A_t) = P[X_{t+1} = x | X_t, A_t]$$

- Rewards are sometimes written as  $R(X_t, A_t, X_{t+1})$ , or just  $R(X_t, A_t)$

- Like with MDPs, the inputs for an MDP are

$\underbrace{S, A}_{\text{state/action spaces}}, \underbrace{T, R}_{\text{transition/reward models}}, \Pi_0$  ← initial state (dist<sup>n</sup> of  $X_0$ )

- To model time varying processes, have  $t$  included in state-space

To model state-dependent action spaces, define  $T$  and  $R$  appropriately

To model 'terminal rewards' in some 'final' state, include dummy actions...

- (Basically, can model everything this way!)

- **Policy**  $\pi = (\pi_1, \pi_2, \dots)$ ,  $\pi_t: S \rightarrow A$  is a collection of mappings (one for each) time from states to actions

## • Optimality Criteria (ie, 'flavors' of MDPs)

MDPs come in different flavors depending on their objective

- **Finite-horizon (Episodic)** - Given known 'horizon'  $H \geq 1$ , for any starting state  $X_0 = x$ , objective is to maximize over all policies  $\pi$ :  
$$V(x) = E_x \left[ \sum_{t=0}^H R_t(X_t, A_t = \pi_t(X_t), U_t) \right]$$

↑ Start at  $X_0 = x$       ← fixed horizon      pick actions from policy  $\pi$

- **Shortest Path problem** - Given (terminal) subset  $U \subset S$ , let  $T_U = \inf \{t \geq 1 \mid X_t \in U\}$ . The objective is to minimize

$$C(x) = E_x \left[ \sum_{t=0}^{T_U-1} R_t(X_t, A_t) \right]$$

↑ 'cost' - sometimes  $R_t(X_t, X_{t+1})$

- **Discounted Reward** - Given discount factor  $\gamma \in (0, 1)$ , objective is

$$V(x) = E_x \left[ \sum_{t=0}^{\infty} \gamma^t R_t(X_t, A_t) \right]$$

Equivalently, given an independent, random horizon  $H \sim \text{Geom}(\gamma)$

$$V(x) = E_x \left[ \sum_{t=0}^H R_t(X_t, A_t) \right]$$

← random horizon  $\sim \text{Geometric}(\gamma)$

- **(Infinite horizon) Average Reward** - Objective is to maximize over all  $(A_t)$

$$V(x) = \limsup_{H \rightarrow \infty} \frac{1}{H} E_x \left[ \sum_{t=0}^H R_t(X_t, A_t) \right]$$

# LP formulations of MDPs

- **Main Idea** - Can 'insulate' future from past decisions by using state-action frequencies as variables

- For finite horizon - Let  $r_t(x,a) \triangleq \mathbb{E}[R_t(X_t=x, A_t=a)]$   
 Consider any policy  $\pi$ , and suppose we 'run' it over many episodes

Now define  $q_t(x,a)$  = fraction of runs which end up in state  $x$  at time  $t$  and policy  $\pi$  plays action  $a$

- Expected reward of  $\pi \equiv V^\pi(x_0) = \sum_{t=0}^H \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} r_t(x,a) q_t(x,a)$

- **Consistency (flow-balance)** -  $q_0(x,a) = 0 \forall x \neq x_0, \sum_{a \in \mathcal{A}} q_0(x_0,a) = 1$

and  $\forall t \geq 1, \forall x \in \mathcal{S}$ :  $\underbrace{\sum_{a \in \mathcal{A}} q_t(x,a)}_{\text{'flow out' of } x, t} = \underbrace{\sum_{x' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} q_{t-1}(x',a') T_{t-1}(x|x',a')}_{\text{'flow in' to } x, t}$

Putting it together we get the LP

Finite-horizon Primal

	$\max \sum_{t=0}^H \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} q_t(x,a) r_t(x,a)$		
dual var	$\text{s.t. } q_0(x,a) = 0 \quad \forall x \neq x_0$		
$V_0(x_0)$	$\sum_{a \in \mathcal{A}} q_0(x_0,a) = 1$		
$V_t(x)$	$\sum_{a \in \mathcal{A}} q_t(x,a) = \sum_{x' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} q_{t-1}(x',a') T_{t-1}(x x',a') \quad \forall x \in \mathcal{S}, t \geq 1$		
	$q_t(x,a) \geq 0 \quad \forall t, x, a$		

We can also now look at the dual LP

$$\begin{aligned} \min \quad & V_0(x_0) \\ \text{s.t. } \quad & \textcircled{*} \quad V_t(x) - \sum_{x' \in S} \tilde{T}_t(x'|x, a) V_{t+1}(x') \geq r_t(x, a) \quad \forall t < H, \forall x, a \\ & V_t(x) \geq 0 \quad \forall t < H, \forall x, a \end{aligned}$$

• Note - If  $X_0 \sim \Pi_0$ , we set  $\sum_a q_0(x, a) = \Pi_0(x) \forall x$  in the primal, and  $\min \sum_x \Pi_0(x) V_0(x)$  as the objective in the dual

• We can simplify  $\textcircled{*}$  in the dual to get

$$V_t(x) \geq \max_{a \in A} \left[ r_t(x, a) + \sum_{y \in S} \tilde{T}_t(y|x, a) V_{t+1}(y) \right] \quad \forall t < H, \forall x \in S$$

Finite-horizon HJB eqns

- This is called the **Bellman optimality condition** (and is a special condition of the more general **Hamilton-Jacobi-Bellman** or **HJB** equation).

- The variables  $\{V_t(x)\}_{x \in S}$  are referred to as the **value function** at  $t$ . Any feasible value fn  $V_t(x)$  induces a corresponding policy  $\Pi_t^V(x) = \underset{a \in A}{\operatorname{argmax}} \left[ r_t(x, a) + \sum_{y \in S} \tilde{T}_t(y|x, a) V_{t+1}(y) \right] \quad \forall t < H, \forall x$

Similarly any policy  $\Pi = (\Pi_t(x))$  induces a corresponding value fn

$$V_t^\Pi(x) = r_t(x, \Pi_t(x)) + \sum_{y \in S} \tilde{T}_t(y|x, \Pi_t(x)) V_{t+1}^\Pi(y) \quad \forall t < H, \forall x$$

(we need as input 'terminal' rewards  $V_H^\Pi(x)$  in both cases)

## LP formulations for other criterion

The advantage of the state-action frequency LP is that it naturally extends to the other flavours of MDPs.

- Discounted rewards
- Consider a **time-invariant** MDP, i.e. with  $R_t = R$  and  $T_t = T$ 
  - Claim - the opt policy can also be taken to be time-invariant
- Suppose we run a policy  $\pi$  over many trials  $j \in \{1, 2, \dots\}$ , where each trial terminates after  $H^j \sim \text{Geom}(\gamma)$  rounds
- As before,  $r(x, a) = \mathbb{E}[R(x, a)]$
- Define  $q(x, a) = \text{avg \# of times action } a \text{ played in state } x$
- Also assume  $X_0 \sim \pi_0$  ← avg over each trial
- Then the MDP  $\equiv$  following LP

$$\begin{array}{ll}
 \max & \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(x, a) q(x, a) \quad \text{Discounted MDP primal} \\
 \text{s.t.} & \\
 \pi_0(x) + \sum_{y \in \mathcal{S}} \sum_{a \in \mathcal{A}} q(y, a) (1-\gamma) T(x|y, a) &= \sum_{a \in \mathcal{A}} q(x, a) \quad \forall x \in \mathcal{X} \\
 \uparrow & \\
 \text{dual or } V(x) & \\
 & q(x, a) \geq 0 \quad \forall x, a
 \end{array}$$

and its dual

$$\begin{array}{ll}
 \min & \sum_{x \in \mathcal{X}} \pi_0(x) V(x) \\
 \text{s.t.} & V(x) \geq r(x, a) + (1-\gamma) \sum_{y \in \mathcal{S}} T(y|x, a) V(y) \quad \forall x, \forall a \\
 & V(x) \geq 0 \quad \forall x
 \end{array}$$

equivalently  $\forall x \in \mathcal{S} \quad V(x) \geq \min_{a \in \mathcal{A}} \left[ r(x, a) + (1-\gamma) \sum_{y \in \mathcal{S}} T(y|x, a) V(y) \right]$

discounted reward HJB eqn