

## Preliminaries

### Initial Notation

- Simplex:  $\Delta_n$  is the simplex of dimension  $n$ .
- States:  $\mathcal{S} = \{s_1, \dots, s_n\}$
- Actions:  $\mathcal{A} = \{a_1, \dots, a_m\}$

- State transitions:

$$G_{t,k}(i, j) = \text{prob}(X_{t+1} | X_t = s_i, U_t = a_k)$$

$$G_{t,k} \geq 0, \quad G_{t,k} \mathbf{1} = \mathbf{1} \quad (\text{row stochastic})$$

- Markovian Policy

$$- \pi_t(s, a) \triangleq \text{prob}(U_t = a | X_t = s)$$

$$- K_t(i, k) \triangleq \pi_t(s_i, a_k), \quad K_t \geq 0, \quad K_t \mathbf{1} = \mathbf{1}$$

$$\pi = (K_0, K_1, \dots) \quad (\text{policy})$$

- Reward and Performance Metric

$$- R_t(s, a) \quad (\text{immediate reward})$$

—

$$v^\pi = \mathbb{E}_{p_o}^\pi \left[ \sum_{t=0}^{N-1} \gamma^t R_t(X_t, U_t) + \gamma^N r_N(x_N) \right] \quad (1)$$

$$= \mathbb{E}_{p_o}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t(X_t, U_t) \right] \quad (2)$$

where  $p_o \in \Delta_n$  is the probability density at  $t = 0$  over states and  $\gamma \in (0, 1]$  is a discount factor.  $\gamma < 1$  when  $N = \infty$ .

### Observations

1.  $p_t \in \Delta_n$  is a prob. dist. over states at time  $t$ .  
Given a decision policy,  $\pi$

$$p_{t+1} = M_{\pi,t}^T p_t \quad (3)$$

where

$$M_{\pi,t} = \sum_{k=1}^m G_{t,k} \odot (K_t e_k \mathbf{1}^T) \quad t = 0, 1, \dots \quad (4)$$

where  $e_k$  is the  $k$ -th standard basis vector and  $\odot$  is the Hadamard product that corresponds to element-wise multiplication. The above propagation

defines a Markov Chain (MC) for the time evolution of the density. Indeed,  $M_{\pi,t}$  is a Markov transition matrix whose element  $[M_{\pi,t}]_{ij}$  is the probability of transitioning from state  $i$  to state  $j$  at time  $t$  under policy  $\pi$ . Note that  $M_{\pi,t}$  is row stochastic.

$$M_{\pi,t} \geq 0, \quad M_{\pi,t} \mathbf{1} = \mathbf{1} \quad p_{t+1}^T = p^T M_{\pi,t}. \quad (5)$$

2. Define the quantities

$$r_{\pi,t} \triangleq (R_t \odot K_t) \mathbf{1} \in \mathbb{R}^n \quad (6)$$

$$v^\pi = \sum_{t=0}^{N \text{ or } \infty} \gamma^t p_t^T r_{\pi,t} \quad (7)$$

where  $p_t$  evolves according to (5)

3. For  $N < \infty$ , finite, letting  $V_{\pi,N} = r_N \in \mathbb{R}^n$ .

$$V_{\pi,t} = r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1}, \quad t = N-1, \dots, 1, 0 \quad (8)$$

We call  $V_{\pi,t} \in \mathbb{R}^n$  the value function at time  $t$ . Note that the value function is initialized as the expected reward at the final time and then propagates backwards in time according to the dynamics. The element  $(V_{\pi,t})_i$  can be thought of as the "reward-to-go" from state  $i$  at time  $t$  when policy  $\pi$  is employed.

4. When a fixed set of actions are given,  $\{U_0, U_1, \dots\}$ , i.e.  $U_t = a_k(t)$ ,  $t = 0, 1, \dots$  it can still be expressed as  $\pi = (K_0, K_1, \dots)$  with  $K_t = \mathbf{1} e_{k(t)}^T$ ,  $t = 0, 1, \dots$

5. In all the following, we'll focus on Markovian policies, i.e. we do *not* consider history dependent ones. Policies can in general be:

- (a) Randomized history dependent (RHD)
- (b) Randomized Markovian (RM)
- (c) Deterministic history dependent (DHD)
- (d) Deterministic Markovian (DM)

Note: Puterman's book has a proof of the fact that MDPs must have DM optimal policies when  $\mathcal{S}$  and  $\mathcal{A}$  are finite sets (see Proposition 4.4.3, p.90?) for MDPs with  $N < \infty$ . This is not necessarily true when the MDP has constraints!

## Finite Horizon MDPs

Problem:  $\max_{\pi} v^\pi$ .

$$v^\pi = p_0^T r_{\pi,0} + \gamma p_1^T r_{\pi,1} + \cdots + \gamma^{N-1} p_{N-1}^T r_{\pi,N-1} + \gamma^N \underbrace{p_N^T r_N}_{p_{N-1}^T M_{\pi,N-1} V_{\pi,N}} \quad (9)$$

$$= \cdots + \gamma^{N-1} p_{N-1}^T \underbrace{\left( r_{\pi,N-1} + \gamma M_{\pi,N-1} V_{\pi,N} \right)}_{V_{\pi,N-1}} \quad (10)$$

$$\begin{array}{cccc} \vdots & \vdots & \vdots & \vdots \\ v^\pi = p_0^T V_{\pi,0} & & & \end{array} \quad (11)$$

Problem:  $\max_\pi p_0^T V_{\pi,0}$

Since  $V_{\pi,N} = r_N$  for all  $\pi$ , we can compute  $V_{\pi,t}$  by backwards induction.

$$V_{\pi,t} = r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1} \quad t = N-1, \dots, 0 \quad (12)$$

Then note that  $M_{\pi,t}$  and  $r_{\pi,t}$  are linear in  $\pi_t$  and

$$e_j^T V_{\pi,t} = e_j^T r_{\pi,t} + \gamma e_j^T M_{\pi,t} \underbrace{V_{\pi,t+1}}_{:=y} \quad (13)$$

where

$$e_j^T r_{\pi,t} = e_j^T (R_t \odot K_t) \mathbf{1} \quad (14)$$

$$= (e_j^T R_t \odot \underbrace{e_j^T K_t}_{j\text{-th row of } K_t}) \mathbf{1} \quad (15)$$

$$e_j^T M_{\pi,t} y = e_j^T \left( \sum_{k=1}^m G_{t,k} \odot (K_t e_k \mathbf{1}^T) \right) y \quad (16)$$

$$= \left( \sum_{k=1}^m e_j^T G_{t,k} \odot \left( \underbrace{e_j^T K_t}_{j\text{-th row of } K_t} e_k \mathbf{1}^T \right) \right) y \quad (17)$$

It follows that  $e_j^T V_{\pi,t}$  linearly depends only on the  $j$ -th row of  $K_t$  (decision variable for policy). Hence in the absence of all other constraints we can solve the following problem to obtain an optimal policy:

For each  $j = 1, \dots, n$ .

$$\max_{e_j^T K_t} e_j^T V_{\pi,t} \quad t = 0, \dots, N-1 \quad (18)$$

$$\text{s.t. } e_j^T K_t \geq 0, \quad e_j^T K_t \mathbf{1} = 1 \quad (19)$$

for each row of  $K_t$ ,  $t = 0, \dots, N-1$ .

Since the cost is linear in the  $j$ -th row one of the optimal choices for  $e_j^T K_t$  is a vertex of the probability simplex in  $\mathbb{R}^m$ !

Note that the overall cost to be maximized is

$$\sum_{t=0}^N \gamma^t p_t^T r_{\pi,t} = p_o^T V_{\pi,0} \quad (20)$$

$$= \sum_{t=0}^{T-1} \gamma^t p_t^T r_{\pi,t} + p_T^T V_{\pi,T} \quad (21)$$

for any  $T = 0, \dots, N-1$ . It follows that

$$\max_{\pi} \sum_{t=0}^N \gamma^t p_t^T r_{\pi,t} = \max_{\substack{K_0, \dots, K_{T-1} \\ K_T, \dots, K_{N-1}}} \sum_{t=0}^{T-1} \gamma^t p_t^T r_{\pi,t} + p_T^T V_{\pi,T} \quad (22)$$

Note  $p_T$  is independent of  $K_T, \dots, K_{N-1}$  and purely depends on  $K_0, \dots, K_{T-1}$ , however  $V_{\pi,T}$  is independent of  $K_0, \dots, K_{T-1}$ .

In this case, meaningful problems are

$$\max_{\pi} \sum_{t=0}^N \gamma^t p_t^T r_{\pi,t} = \max_{\pi} p_0^T V_{\pi,0} \quad \text{for a given } p_0. \quad (23)$$

OR

$$\max_{\pi} \min_{p_0} \sum_{t=0}^N \gamma^t p_t^T r_{\pi,t} = \max_{\pi} \min_{p_0} p_0^T V_{\pi,0} \quad (24)$$

Since  $e_j^T V_{\pi,t}$  linearly depends on  $e_j^T K_t$  only, it can be maximized with it's choice. Also, since

$$e_j^T V_{\pi,t} = e_j^T \underbrace{r_{\pi,t}}_{\substack{\text{linear} \\ \text{in } e_j^T K_t}} + \gamma e_j^T \underbrace{M_{\pi,t}}_{\substack{\text{linear} \\ \text{in } e_j^T K_t}} V_{\pi,t+1} \quad (25)$$

No matter how  $K_t$  is chosen to maximize  $e_j^T V_{\pi,t}$ , we have to maximize each component of  $V_{\pi,t+1}$  separately simply since it's possible as above. More precisely,

$$\pi^* = \arg \max_{\pi} \min_{p_0 \in \Delta_n} p_0^T V_{\pi,0} \quad (26)$$

$$= \arg \max_{\pi} e_j^T V_{\pi,0} \quad j = 1, \dots, n \quad (27)$$

since  $e_j \in \Delta_n$  for  $j = 1, \dots, n$ . It follows that

$$V_{\pi^*,0} \geq V_{\pi,0} \quad \forall \pi \quad (28)$$

$$\Rightarrow p_0^T V_{\pi^*,0} \geq p_0^T V_{\pi,0} \quad \forall \pi \quad (29)$$

$$\Rightarrow \pi^* = \max_{\pi} p_0^T V_{\pi,0} \quad \forall p_0 \in \Delta_n \quad (30)$$

It follows that

$$\max_{\pi} \min_{p_0 \in \Delta} p_0^T V_{\pi,0} \quad (31)$$

is a proper problem!  
Now, recall that

$$V_{\pi,N} = r_N \quad (32)$$

$$V_{\pi,t} = r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1} \quad (33)$$

Since, as shown above, we have that

$$V_{\pi,t} = \begin{bmatrix} \max_{\pi} e_1^T V_{\pi,t} \\ \vdots \\ \max_{\pi} e_n^T V_{\pi,t} \end{bmatrix} \quad (34)$$

for a given  $V_{\pi,t+1}$  we have that

$$V_{\pi^*,t} \geq V_{\pi,t} \quad \forall \pi, \quad t = 0, \dots, N-1 \quad (35)$$

The conclusion is that no matter how  $p_t$  evolves in time, the optimal policy can be computed in one-shot via dynamic programming (DP).

### Dynamic Programming for Finite-Horizon MDPs

- Initialize.  $V_N^* = r_N$
- For  $t = N-1, \dots, 0$ , iteratively compute

$$e_j^T K_t^* = \arg \max_{e_j^T K_t^*, e_j \in \Delta_n} e_j^T r_{\pi,t} + e_j^T \gamma M_{\pi,t} V_{\pi,t+1}^* \quad (36)$$

$$V_{\pi,t}^* = \max_{e_j^T K_t^*, e_j \in \Delta_n} e_j^T r_{\pi,t} + e_j^T \gamma M_{\pi,t} V_{\pi,t+1}^* \quad (37)$$

Since the above maximization is an LP over  $\Delta_n$ , we can always find an optimal policy that is deterministic Markovian (DMP) in the absence of constraints.

## Infinite Horizon MDPs

In the infinite horizon case, we assume the transition kernel and rewards are constant over time, i.e.  $G_t = G$  and  $R_t = R$  for all  $t$ .

### Theorem 1 (Existence and Uniqueness)

*The following equation has a unique solution*

$$T(V) = \max_{\pi} (r_{\pi} + \gamma M_{\pi} V) \quad (38)$$

**Proof 1** The proof proceeds by using the Banach fixed point theorem on the Bellman operator  $T(V)$ . Consider  $V_1 \rightarrow V_1^*$  and  $V_2 \rightarrow V_2^*$  for any  $V_1, V_2 \in \mathbb{R}^n$

$$V_k^* = \max_{\pi} (r_{\pi} + \gamma M_{\pi} V_k), \quad k = 1, 2 \quad (39)$$

We consider the difference  $\|T(V_1) - T(V_2)\|_{\infty}$  under the infinity norm.

$$\|T(V_2) - T(V_1)\|_{\infty} = \left\| \max_{\pi} (r_{\pi} + \gamma M_{\pi} V_2) - \max_{\pi} (r_{\pi} + \gamma M_{\pi} V_1) \right\|_{\infty} \quad (40)$$

$$= \max_{j=1 \dots n} \left| \max_{\pi} e_j^T (r_{\pi} + \gamma M_{\pi} V_2) - \max_{\pi} e_j^T (r_{\pi} + \gamma M_{\pi} V_1) \right| \quad (41)$$

$$\leq \max_j \max_{\pi} |e_j^T (r_{\pi} + \gamma M_{\pi} V_2) - e_j^T (r_{\pi} + \gamma M_{\pi} V_1)| \quad (42)$$

$$\leq \max_j \max_{\pi} \gamma |e_j^T M_{\pi} V_2 - e_j^T M_{\pi} V_1| \quad (43)$$

$$= \max_{\pi} \gamma \|M_{\pi} (V_2 - V_1)\|_{\infty} \quad (44)$$

$$= \max_{\pi} \gamma \|V_2 - V_1\|_{\infty} \quad (45)$$

where we have used two facts:

$$\max_x f(x) - \max_y g(y) \leq \max_x f(x) - g(x) \quad (46)$$

and

$$\|Mv\|_{\infty} \leq \|M\|_{\infty} \|v\|_{\infty} = \|v\|_{\infty} \quad (47)$$

since  $\|M\|_{\infty} = \max_j \|e_j^T M\|_1 = 1$ . Thus we have that  $T(V)$  is a contractive mapping for  $\gamma \in [0, 1)$ . By the Banach fixed point theorem, it follows that  $V = T(V)$  has a unique solution  $V^*$  and that the sequence  $V_{k+1} = T(V_k)$  for any  $V_0$  will converge to  $V^*$ .

Explicitly we have that

$$\|V_{k+1} - V_k\| \leq \gamma \|V_k - V_{k-1}\| \quad (48)$$

$$\Rightarrow \|V_k - V^*\| \leq \gamma^k \|V_0 - V^*\| \quad (49)$$

This proof leads to a straight-forward technique for computing the value function called *value iteration*.

### Value Iteration

- Pick  $V_0$ .
- Compute  $V_{k+1} = T(V_k)$ ,  $k = 0, 1, \dots$
- Stop when  $\|V_{k+1} - V_k\|$  is within some desired tolerance.

Properties of  $V^*$  and the corresponding  $\pi^*$ :

$$v^\pi = \sum_{t=0}^{\infty} \gamma^t \underbrace{p_t^T}_{\text{linear in } \pi} \underbrace{r_{\pi,t}}_{\text{linear in } \pi} \quad (50)$$

For any policy, since  $\exists \alpha_1, \alpha_2 > 0$  such that

$$\|r_{\pi,t}\| \leq \alpha_1, \quad \|p_t\| \leq \alpha_2 \quad \Rightarrow \quad |p_t^T r_{\pi,t}| \leq \underbrace{\alpha_1 \alpha_2}_{:=\alpha} \quad (51)$$

$$\Rightarrow \quad \sum_{t=0}^{\infty} \gamma^t |p_t^T r_{\pi,t}| = \alpha \sum_{t=0}^{\infty} \gamma^t = \alpha \frac{1}{1-\gamma} \quad (52)$$

It follows that  $v^\pi$  is absolutely convergent and thus  $v^\pi = c$  for some  $c$ , i.e. all policies have a finite reward for any  $p_0 \in \Delta_n$ .

$$p_t^T = p_0^T M_{\pi,0} \cdots M_{\pi,t-1} \quad (53)$$

$$v^\pi = \sum_{t=0}^{\infty} \gamma^t p_0^T \left( \underbrace{M_{\pi,0} \cdots M_{\pi,t-1}}_{:=\tilde{M}_{\pi,t-1}} \right) r_{\pi,t} \quad (54)$$

$$= p_0^T \left[ r_{\pi,0} + \gamma M_{\pi,0} r_{\pi,1} + \gamma^2 M_{\pi,0} M_{\pi,1} r_{\pi,2} + \cdots \right] \quad (55)$$

$$= p_0^T \left[ r_{\pi,0} + \gamma M_{\pi,0} \left( \underbrace{r_{\pi,1} + \gamma M_{\pi,1} r_{\pi,2} + \cdots}_{:=V_{\pi,1}} \right) \right] \quad (56)$$

$$= p_0^T \left[ \underbrace{r_{\pi,0} + \gamma M_{\pi,0} V_{\pi,1}}_{:=V_{\pi,0}} \right] \quad (57)$$

where

$$V_{\pi,1} = r_{\pi,1} + \gamma M_{\pi,1} \left( \underbrace{r_{\pi,2} + \gamma M_{\pi,2} r_{\pi,3} + \cdots}_{:=V_{\pi,2}} \right) \quad (58)$$

$$\Rightarrow \quad V_{\pi,t} \triangleq r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1} \quad (59)$$

Note  $V_{\pi,t}$  are well-defined for all  $\gamma \in [0,1)$  due to absolute convergence. It follows that  $v^\pi = p_0^T V_{\pi,0}$ . As before  $\max_\pi v^\pi = \max_\pi e_j^T V_{\pi,0}$   $j = 1, \dots, n$ .

$$V_{\pi,0} = r_{\pi,0} + \gamma M_{\pi,0} V_{\pi,1} \quad (60)$$

$$\vdots \quad (61)$$

$$V_{\pi,t} = r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1} \quad t = 0, 1, 2, \dots \quad (62)$$

Note: Consider a stationary policy,  $\pi = (\pi, \pi, \dots)$  with abuse of notation

$$V_\pi = r_\pi + \gamma M_\pi V_\pi \quad (63)$$

$$\Rightarrow V_\pi = (I - \gamma M_\pi)^{-1} r_\pi \quad (64)$$

The spectrum of  $M_\pi$  is in the unit circle  $\Rightarrow I - \gamma M_\pi$  cannot have a zero eigenvalue. Thus  $I - \gamma M_\pi$  is invertible and for stationary policies for stationary processes, we have that

$$V_\pi = (I - \gamma M_\pi)^{-1} r_\pi \quad (65)$$

$$= \sum_{t=0}^{\infty} \gamma^t M_\pi^t r_\pi \quad (66)$$

Since, for Markovian policies,

$$V_{\pi,t} = r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1} \quad (67)$$

For optimal, Markovian policies

$$V_{\pi,t}^* = \max_{\pi} r_{\pi,t} + \gamma M_{\pi,t} V_{\pi,t+1}^* \quad (68)$$

i.e.  $V_{\pi,t}^* = T(V_{\pi,t+1}^*)$ .