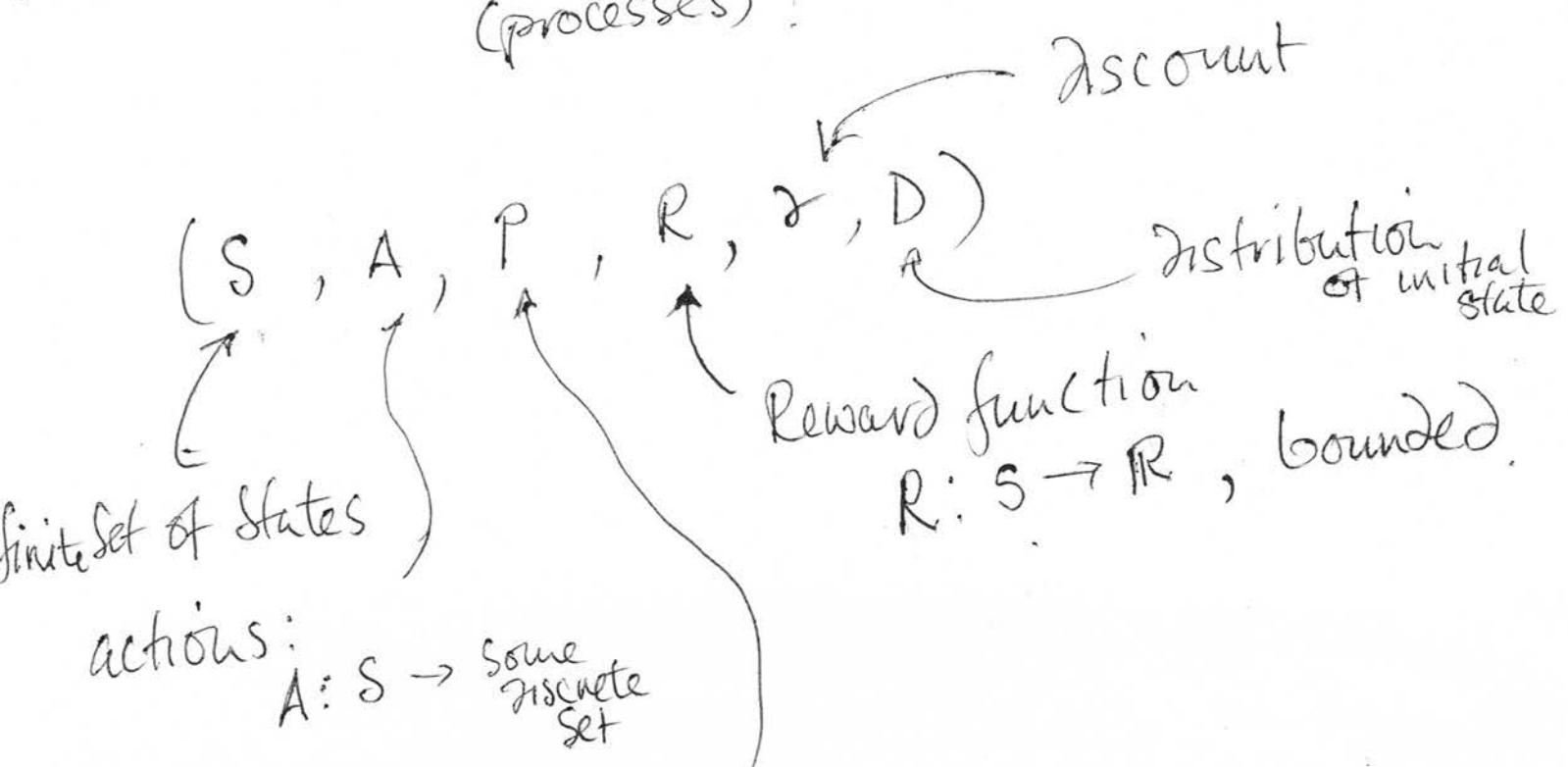


# Markov Decision problems: (processes)!



model  $P(S_{i+1} = s \mid S_i = \hat{s}, a(s_i) = \hat{a})$

## Policy

$$\pi : S \rightarrow A$$

(i.e. in state  $S$ , do this)

General, very useful model for action + control

- + Can ~~be~~ extend.
- + Note moral similarity to LQR.

(Where we ~~could~~ had states actions = controls

But transitions were deterministic

Model

- Start in  $s_0 \sim D, [total\ reward] = 0$

→ - choose  $a_i$

-  $s_{i+1} \sim P(s_{i+1} | s_i, a_i)$

- get reward  $R(s_{i+1})$

~~$[total\ reward] \rightarrow \gamma [total\ reward] + R(s_{i+1})$~~

Notice effect of discount — long ago rewards have less weight.

want to max

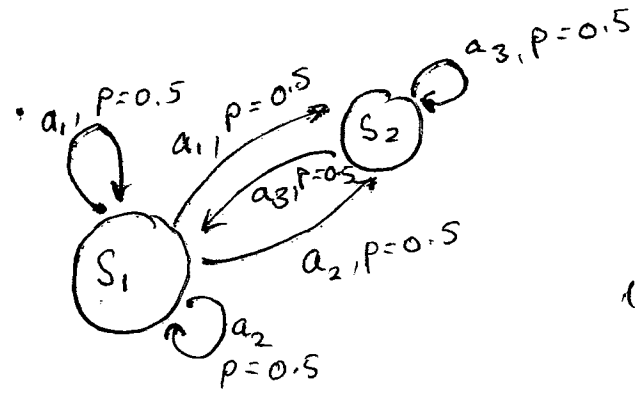
$$\sum_{t=0}^{\infty} \gamma^t R(s_t)$$

↑ discounted total reward

# Natural Question:

- What is the best policy  $\pi^*$  for a given problem?  
(This might not be unique)

eg



is as good as  $\pi(s_1) = a_1$   
 $\pi(s_1) = a_2$

# Subordinate Question:

assume policy is  $\pi$ ; what is

Value of states?

$$i.e. V(s) = E \left[ \text{Reward to be obtained by following } \pi \text{ from } s \right]$$

where expectation is over  $P$ .

Value is important, because its linked to policy. imagine we know  $V_{\pi^*}(s)$ .

then we can recover  $\pi^*$  choose action with best expected value of  $V_{\pi^*}(s)$ .

Notice  $a_i$  has expected value  $\sum_s P(s|\hat{s}, a_i) V_{\pi^*}(s)$

Notice:  $V_{\pi^*}(s)$  is usually more than  $R(s)$  because you get rewards from future actions

~~But  $\gamma = 0 \implies V_{\pi^*}(s) = R(s)$~~

get  $V_{\pi^*}(s) = R(s) + \gamma \max_a \sum P(s'|s, a) V_{\pi^*}(s')$

AND  $\pi^*(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) V_{\pi^*}(s')$

Bellman equation

# Estimating the optimal value fn.

## Value Iteration:

init  $V^{(0)}(s) = 0$

update  $V^{(n+1)}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{(n)}(s)$

Then  $V^{(k)} \rightarrow V^{\pi^*}$  as  $k \rightarrow \infty$

## Proof:

write  $B[V] = R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V(s)$   
↑  
Bellman Backup operator.

lemma below shows

$$\max_{\theta} |B[V_1](s) - B[V_2](s)| \leq \gamma \max_s |V_1(s) - V_2(s)|$$

Note that

$$B[V_{\pi^*}] = V_{\pi^*}$$

assume lemma

$$\begin{aligned} \text{then } & |B[V^{(u)}] - B[V_{\pi^*}]|_{\infty} \\ &= |B[V^{(u)}] - V_{\pi^*}|_{\infty} \\ &\leq \gamma |V^{(u)} - V_{\pi^*}|_{\infty} \end{aligned}$$

and  $0 \leq \gamma < 1$

so  $V^{(u+1)} = B[V^{(u)}]$  is closer to  $V_{\pi^*}$  than  $V^{(u)}$   $\square$

lemma:

$$|B[V_1] - B[V_2]|_\infty = \gamma \left| \left[ \max_a \sum_{s'} P(s'|s,a) V_1(s') \right] - \left[ \max_a \sum_{s'} P(s'|s,a) V_2(s') \right] \right|_\infty$$

$$\leq \gamma \max_a \left| \sum_{s'} P(s'|s,a) V_1(s') - \sum_{s'} P(s'|s,a) V_2(s') \right|_\infty = A$$

because  $\left| \max_x f(x) - \max_x g(x) \right| \leq \max_x |f(x) - g(x)|$

$$A = \gamma \max_a \sum_{s'} P(s'|s,a) |V_1(s') - V_2(s')|$$

$$\leq \gamma \max_s |V_1(s) - V_2(s)|$$

(P is non neg, and sums to 1)

# Convergence :

assume rewards bounded by  $R_{max}$

$$V^*(s) \leq \sum_{t=0}^{\infty} \gamma^t R_{max} = \frac{R_{max}}{1-\gamma}$$

↑

so  $\max_{s \in S} |V^k(s) - V^*(s)| \leq \frac{\gamma^k R_{max}}{1-\gamma}$

(i.e. linear ~~fast~~; faster for small  $\gamma$ )

# Asynchronous Version :

- notice what is above assumes all states are updated together

- i.e. compute  $\tilde{V}(s) = B[V^{(n)}]$   
 then  $V^{(n+1)} = \tilde{V}(s)$

- but you could visit states in turn, updating each

- Converges
- usually better



# Policy Iteration:

- initialize w/ policy, perhaps random  $\pi^{(0)}$
- Compute  $V_{\pi}$

• Notice this involves solving a linear system

$$V_{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V_{\pi}(s')$$
$$V_{\pi} = R + \gamma T V_{\pi}$$
$$(I - \gamma T) V_{\pi} = R$$
$$V_{\pi} = [I + \gamma T + \gamma^2 T^2 + \gamma^3 T^3 \dots] R$$

↳ Neumann Series

- update  $\pi$  to be greedy w.r.t  $V_{\pi}$   
i.e.  $\pi(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) V_{\pi}(s')$
- Iterate

Policy iteration converges

- to exact optimal policy
- eventually.

In practice, policy iteration is usually preferred, particularly if transition probs are sparse

MDP's as LP's :

$$V(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s'|s, a) V(s')$$

implies ~~equation~~ to

---


$$V(s) \geq R(s) + \gamma \sum_{s' \in S} P(s'|s, a) V(s')$$

↑ one constraint for each A

So we can seek

$$\begin{aligned} \min & \sum_s V(s) \\ \text{s.t.} & V(s) \geq R(s) + \gamma \sum_{s'} P(s'|s,a) V(s') \end{aligned} \quad \left| \leftarrow \text{one for each } s, a \right.$$

and this is the optimal value fn.

Write  $\mathbf{1}^T V$

$$\text{s.t.} \quad (\mathbf{I} - T_a) V - R \geq 0 \quad a = 1 \dots \#A$$

V is #S vector  
one  $T_a$  per a.

Now consider the dual

$\lambda_a$   
↑  
LM, #S dimensional, 1 per a

$$\begin{aligned} \max & \sum_a \lambda_a^T R \\ \text{s.t.} & \sum_a \lambda_a^T (\mathbf{I} - T_a)^T + \mathbf{1} = 0 \\ & \lambda_a \geq 0 \end{aligned}$$

What are Dual vars?

1) Notice one useful interpretation of  $\gamma$ .

$$(1 - \gamma) = P[\text{MDP stops}]$$

$$\begin{aligned}
E[\text{reward}] &= \text{reward at start} \\
&= \text{reward at 1st state} P(1 \text{ step}) \\
&+ \text{reward at 2nd state} P(2) \\
&+ \dots \\
&= R(s_0) + R(s_1)\gamma + R(s_2)\gamma^2 + \dots
\end{aligned}$$

2) write Primal objective

$$\frac{1^T V}{|S|} \leftarrow \# \text{ states.}$$

Dual becomes

$$\max \sum_a \lambda_a^T R$$

$$\text{s.t. } \sum_a (1 - P_a)^T \lambda_a + \frac{1}{|S|} = 0$$

$$\lambda_a \geq 0$$

Now consider the expected number of times we are in  $S$ , take a  
 $= N(s, a)$ .

~~$N(s, a) = \frac{1}{|S|}$~~

expected # of times we ~~are in~~ arrive in  $S$  is

$$\frac{1}{|S|} + \gamma \sum_{s'} \sum_a P(s'|a, s) N(s', a)$$

Because we have a uniform prob first state is  $S$

joint Prob  $P(s'|s, a)$

Notice :  $\sum_a N(s, a) =$  expected number of times we leave from  $S$

BUT  $\sum_a N(s, a) = \frac{1}{|S|} + \gamma \sum_{s'} \sum_a P(s'|s, a) N(s', a)$

these are the constraints on dual Vars!

So

$$\lambda_a = \lambda_a(s) = N(s, a).$$

These encode the optimal policy:

$$\begin{aligned} \pi_*(s) &= \text{take the most common action at } s \\ &= \underset{a}{\operatorname{argmax}} N(s, a). \end{aligned}$$

Natural follow up questions

- What happens if you don't know  $P, R$ ?

(Reinforcement learning:

- roughly: act, see what happens  
keep notes

- What happens if you don't know  $P, R$  BUT see a skilled actor?

(our problem)