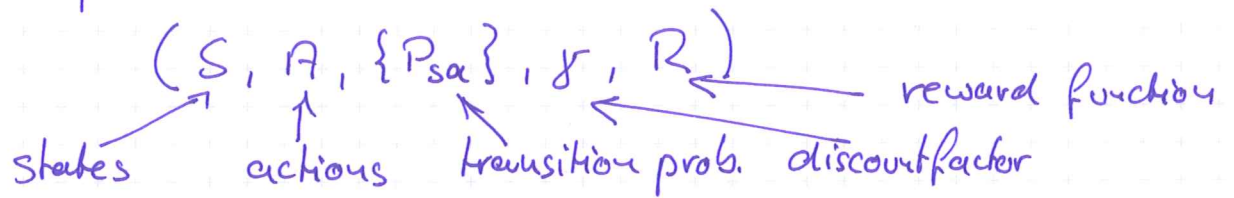


## 7.2. Reinforcement Learning

### 7.2.1. Markov Decision Processes (MDP)

Recapitulation MDP:

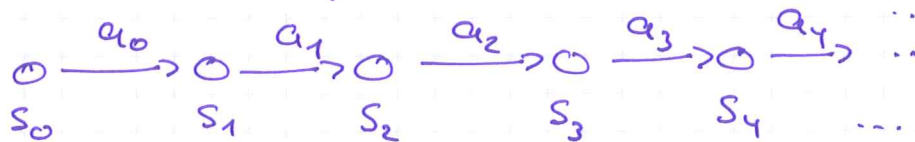


Payoff:  $R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$

Objective:

$$\max_{\text{actions} \in A} E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

Dynamics of the process:



Note that  $s_0, s_1, s_2, \dots$  are random variables.

Policy is a function

$$\pi : S \rightarrow A : s \mapsto \pi(s)$$

In state  $s \in S$  take action  $a = \pi(s)$

Value function for a policy  $\pi$

$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi]$$

Expected total payoff upon starting in  $s$ , applying policy  $\pi$ .

For any policy  $\pi$  the Bellman equations hold:

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s, \pi(s)}(s') V^\pi(s'), \quad s \in S \quad (\text{BE})$$

RHS immediate reward  $R(s)$  + expected sum of future discounted rewards.

Solving for  $V^\pi$  for fixed policy  $\pi$  means to solve  $|S|$  linear equations with  $|S|$  variables provided  $|S| < \infty$ ,  $S = \{1, \dots, m\}$ .

Write  $\underline{R} = (R(s_1), \dots, R(s_m))^T \in \mathbb{R}^m$

$$\Psi^\pi = \begin{pmatrix} P_{1, \pi(s_1)} \\ \vdots \\ P_{m, \pi(s_m)} \end{pmatrix} \in \mathbb{R}^{m \times m}$$

$$\underline{V}^\pi = (V^\pi(s_1), \dots, V^\pi(s_m))^T \in \mathbb{R}^m \quad (\text{variables})$$

Then (BE) in matrix form

$$\underline{V}^{\pi} = \underline{R} + \gamma \Psi^{\pi} \underline{V}^{\pi}$$

$$\Leftrightarrow \underline{V}^{\pi} = (\underline{I}_m - \gamma \Psi^{\pi})^{-1} \underline{R}$$

provided the inverse exists.

Optimal value function:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$= \max_{\pi} E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi]$$

Max. expected payoff over all policies  $\pi$  starting in  $s \in S$ .

It holds, analogous to (BE)

$$V^*(s) = R(s) + \max_{a \in A} \sum_{s' \in S} P_{s,a}(s') V^*(s') \quad (*)$$

Define a partial ordering over policies.

$$\pi \succeq \pi' \text{ if } V^{\pi}(s) \geq V^{\pi'}(s) \text{ for all } s \in S.$$

Theorem 7.1

a) There exists an optimal policy  $\pi^*$  with

$$\pi^* \succeq \pi \text{ for all policies } \pi.$$

b) All optimal policies achieve the same value

$$V^{\pi^*}(s) = V^*(s).$$

⊥

## 7-2-2. Computing the optimum policy

Algorithm policy iteration (PI) to find the opt.  $\pi^*$ :

- (PI)
1. Initialize  $\pi$  randomly
  2. Repeat until convergence {
    - (a)  $\underline{V} := \underline{V}^\pi$
    - (b) For each  $s \in S$ , let  ~~$\pi(s)$~~ 

$$\pi(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sA}(s') V(s')$$

(a) For fixed policy  $\pi$ ,  $\underline{V}^\pi$  can be computed from (BE), see 7.2.1.

(b) Note that  $V(s)$  is w.r.t. to actual policy  $\pi$ .  
Update of  $\pi$  is greedy w.r.t.  $\underline{V}$ .

After a finite number of iterations  $\underline{V}$  converges to  $\underline{V}^*$  and  $\pi$  to  $\pi^*$ .



### 7.2.3 Model Learning for MDP

In practice trans. prob.  $P_{s,a}$  and sometimes the reward fun.  $R$  are unknown.

Solution: estimate  $P_{s,a}(s')$ ,  $s' \in S$ , and  $R(s')$  from data.

Carry out trials/experiments/simulations:

$$s_0^{(l)} \xrightarrow{a_0^{(l)}} s_1^{(l)} \xrightarrow{a_1^{(l)}} s_2^{(l)} \xrightarrow{a_2^{(l)}} \dots, \quad l=1,2,\dots$$

$l$  counts the no. of trials.

Estimate 
$$P_{s,a}(s') = \frac{\# \text{ times took action } a \text{ in state } s \text{ and got to } s'}{\# \text{ times took action } a \text{ in state } s}$$

In case  $\frac{0}{0}$  set  $P_{s,a}(s') = \frac{1}{|S|}$  (uniform distribution)

Similarly  $R(s)$  can be estimated from observed data.

Possible algorithm for learning a MDP with unknown  $P_{s,a}$

1. Initialize  $\pi$  randomly
2. Repeat {
  - (a) Execute  $\pi$  in the MDP for some no. of trials
  - (b) Upgrade estimates  $P_{s,a}$  (and potentially  $R$ )
  - (c) Update  $\pi$  by algorithm (PI)