

## Lecture - 9

### \* Policy evaluation

- Given a policy  $\pi$ , find  $V^\pi$
- $P$  &  $R$  are known

Method:

- ① Solve the system of linear eq<sup>n</sup> produced by the Bellman equation

$$V^\pi(s) = \sum_a \pi(s,a) \sum_{s'} P(s,a,s') \left[ R(s,a,s') + \gamma V^\pi(s') \right]$$

### ② Dynamic Programming

- sequence of approximations of

$V^\pi$   
e.g.  $V_0^\pi, V_1^\pi, V_2^\pi, \dots$

Chosen

arbitrarily

except  $V_0^\pi(s) = 0$

i.e. for terminal

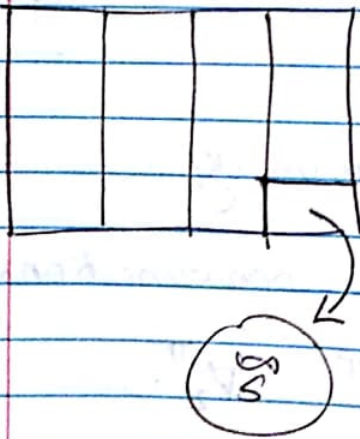
# Dynamic programming . . . . .

$$\textcircled{3} \hat{V}_{k+1}^{\pi}(s) = \sum_a \pi(s,a) \sum_{s'} P(s,a,s') \left[ R(s,a,s') + \gamma \hat{V}_k^{\pi}(s') \right]$$

Properties -

- $\hat{V}_k^{\pi} = V^{\pi}$  is a fixed-point.
- $\hat{V}_k^{\pi} \rightarrow V^{\pi}$  as  $k \rightarrow \infty$  for finite MDPs
- One pass over state space is a "full backup". A single state update is called "backup".

Example :-



$R_t = -1$  always  
 $\gamma = 1$

$$\pi = \begin{bmatrix} \downarrow & \downarrow & \downarrow & \downarrow \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \rightarrow & \rightarrow & \rightarrow & \rightarrow \end{bmatrix}$$

$$\hat{V}_0^{\pi} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \Rightarrow \hat{V}_1^{\pi} = \begin{bmatrix} & & & s' \\ & & & s'' \\ & & & s''' \\ 0 & & & s'''' \end{bmatrix}$$

$$V_1^\pi(s) = s'(-1 + \gamma \hat{V}_k^\pi(s)) \rightarrow \text{Prob.} = \text{zero}$$

$$+ \dots + \underbrace{(-1 + \gamma \hat{V}_k^\pi(s'''))}_{0} \rightarrow \text{Prob.} = 1$$

$$\therefore \hat{V}_1^\pi = \begin{bmatrix} -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & 0 \end{bmatrix}$$

$$\text{Now, } \hat{V}_2^\pi = \begin{bmatrix} -2 & -2 & -2 & -2 \\ -2 & -2 & -2 & -2 \\ -2 & -2 & -2 & -1 \\ -2 & -2 & -1 & 0 \end{bmatrix}$$

$$\dots \hat{V}_6^\pi = \begin{bmatrix} -6 & -5 & -4 & -3 \\ -5 & -4 & -3 & -2 \\ -4 & -3 & -2 & -1 \\ -3 & -2 & -1 & 0 \end{bmatrix}$$

\* In-place implementation

- when updating  $\hat{V}_k^\pi(s)$  store it back in same table
- update states in any order
- Asynchronous updates

\* Policy improvement

$$\hat{q}_{k+1}^{\pi}(s,a) = \sum P(s,a,s') [R(s,a,s') + \gamma \pi(s',a') \cdot \hat{q}_k^{\pi}(s',a')]$$

Let  $\pi$  &  $\pi'$  be deterministic policies such that for all states

$$q^{\pi'}(s, \pi'(s)) \geq v^{\pi}(s) = q^{\pi}(s, \pi(s))$$

Then  $\pi' \geq \pi$

$$\forall s \in \mathcal{S}, v^{\pi'}(s) \geq v^{\pi}(s)$$