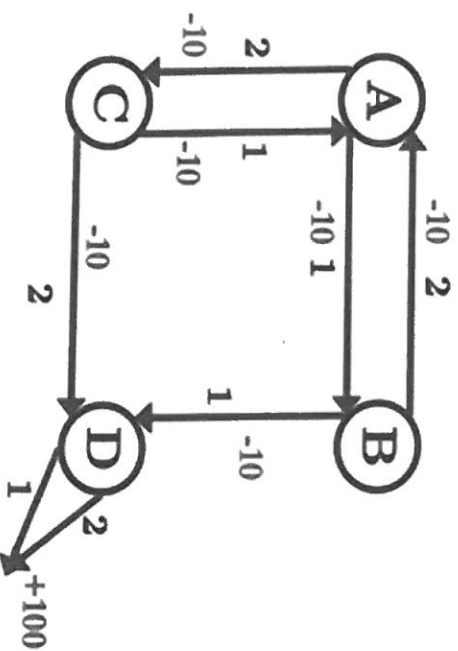


## EXAMPLES



Goal = Terminal state

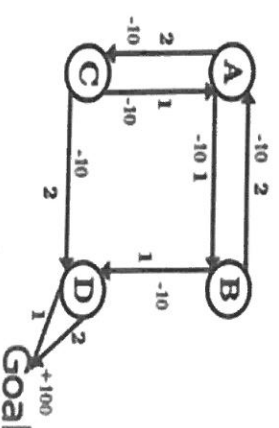
- 4 states
- 2 possible actions in each state. [E.g in A: 1) go to B or 2) go to C ]
- $P(s' | s, a) = (0.9, 0.1)$  with 10% we go to a wrong direction

Rewards associated with each state-action pair are given on the arcs.

Calculating state values for a policy  $\pi_1$ :

$\pi_1$  : always choosing Action 1

$$\begin{aligned} \pi(1|A) &= 1 & \pi(2|A) &= 0 \\ \pi(1|B) &= 1 & \pi(2|B) &= 0 \\ \pi(1|C) &= 1 & \pi(2|C) &= 0 \\ \pi(1|D) &= 1 & \pi(2|D) &= 0 \end{aligned}$$



$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a) + \gamma V^\pi(s')]$$

$$V^\pi(D) = 100$$

$$V^\pi(B) = 1 \cdot 0.9 \cdot [-10 + V^\pi(D)] + 1 \cdot 0.1 \cdot [-10 + V^\pi(A)]$$

$$V^\pi(C) = 1 \cdot 0.9 \cdot [-10 + V^\pi(A)] + 1 \cdot 0.1 \cdot [-10 + V^\pi(D)]$$

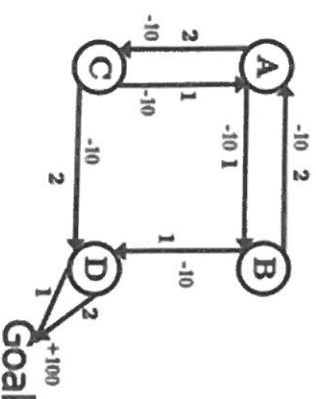
$$V^\pi(A) = 1 \cdot 0.9 \cdot [-10 + V^\pi(B)] + 1 \cdot 0.1 \cdot [-10 + V^\pi(C)]$$

$$\begin{aligned} V^\pi(A) &= 75.61 \\ V^\pi(B) &= 87.56 \\ V^\pi(C) &= 68.05 \\ V^\pi(D) &= 100 \end{aligned}$$

Calculating state values for a policy  $\pi_2$ :

$\pi_2$  : always choosing Action 2

$$\begin{aligned} \pi(1|A) &= 0 & \pi(2|A) &= 1 \\ \pi(1|B) &= 0 & \pi(2|B) &= 1 \\ \pi(1|C) &= 0 & \pi(2|C) &= 1 \\ \pi(1|D) &= 0 & \pi(2|D) &= 1 \end{aligned}$$



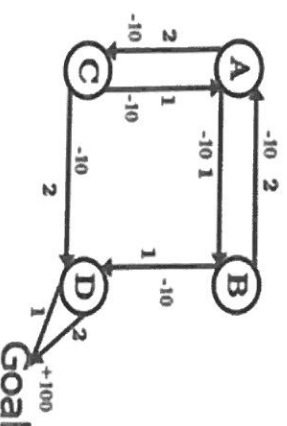
Similarly as before:

$$\begin{aligned} V^{\pi_2}(A) &= 75.61 \\ V^{\pi_2}(B) &= 68.05 \\ V^{\pi_2}(C) &= 87.56 \\ V^{\pi_2}(D) &= 100 \end{aligned}$$

Calculating state values for a policy  $\pi_3$ :

$\pi_3$  : mixed

$$\begin{aligned} \pi(1|A) &= 0.4 & \pi(2|A) &= 0.6 \\ \pi(1|B) &= 1 & \pi(2|B) &= 0 \\ \pi(1|C) &= 0 & \pi(2|C) &= 1 \\ \pi(1|D) &= 1 & \pi(2|D) &= 0 \end{aligned}$$



$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a) + \gamma V^\pi(s')] ]$$

$$\begin{aligned} V^\pi(D) &= \frac{1}{100} \\ V^\pi(B) &= 1 \cdot 0.9 \cdot [-10 + V^\pi(D)] \\ &\quad + 1 \cdot 0.1 \cdot [-10 + V^\pi(A)] \\ V^\pi(C) &= 1 \cdot 0.9 \cdot [-10 + V^\pi(D)] \\ &\quad + 1 \cdot 0.1 \cdot [-10 + V^\pi(A)] \\ V^\pi(A) &= 0.4 \cdot 0.9 \cdot [-10 + V^\pi(B)] \\ &\quad + 0.4 \cdot 0.1 \cdot [-10 + V^\pi(C)] \\ &\quad + 0.6 \cdot 0.9 \cdot [-10 + V^\pi(C)] \\ &\quad + 0.6 \cdot 0.1 \cdot [-10 + V^\pi(B)] \end{aligned}$$

⇔

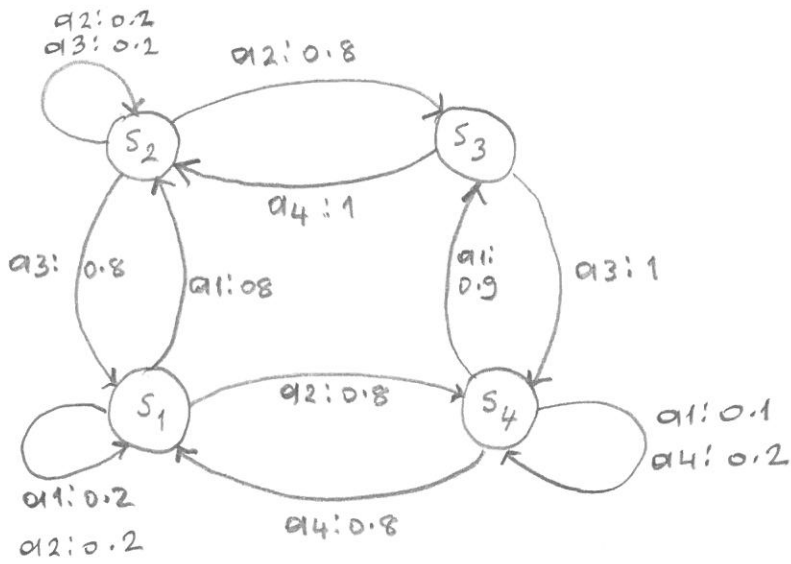
$$\begin{aligned} V^{\pi_3}(A) &= 77.78 \\ V^{\pi_3}(B) &= 87.78 \\ V^{\pi_3}(C) &= 87.78 \\ V^{\pi_3}(D) &= 100 \end{aligned}$$

## Comparing the 3 policies

	$\pi_1$	$\pi_2$	$\pi_3$
A	75.61	75.61	77.78
B	87.56	68.05	87.78
C	68.05	87.56	87.78
D	100	100	100

$\pi_1 \leq \pi_3$  and  $\pi_2 \leq \pi_3 \Rightarrow \pi_3$  is optimal among these 3 policies  
 $\pi_1$  and  $\pi_2$  are not comparable

Q.6. Consider The following MDP



S	r(S)
S <sub>1</sub>	0
S <sub>2</sub>	0
S <sub>3</sub>	1
S <sub>4</sub>	0

$\gamma = 0.5$

t	0	1	2	3	4
V(S <sub>1</sub> )	0	0	0	0.18	0.206
V(S <sub>2</sub> )	0	0	0.4	0.44	0.52
V(S <sub>3</sub> )	0	1	1	1.225	1.2
V(S <sub>4</sub> )	0	0	0.45	0.473	0.57

$$V^*(S_1) = 0 + \gamma \cdot \max \left\{ 0.2 \times V^*(S_1) + 0.8 \times V^*(S_2); 0.2 \times V^*(S_1) + 0.8 \times V^*(S_4) \right\}$$

$$V^*(S_2) = 0 + \gamma \cdot \max \left\{ 0.2 \times V^*(S_2) + 0.8 \times V^*(S_3); 0.2 \times V^*(S_2) + 0.8 \times V^*(S_1) \right\}$$

$$V^*(S_3) = 1 + \gamma \cdot \max \left\{ V^*(S_2), V^*(S_4) \right\}$$

$$V^*(S_4) = 0 + \gamma \cdot \max \left\{ 0.1 \times V^*(S_4) + 0.9 \times V^*(S_3); 0.2 \times V^*(S_4) + 0.8 \times V^*(S_1) \right\}$$