

MDP

① Bellman-EQ

$$\begin{aligned}
 * U(s) &= \max_{a \in A(s)} \sum_{s'} P(s'/s, a) [R(s, a, s') + \gamma U(s')] \\
 &= \max_{a \in A(s)} R(s, a) + \sum_{s'} P(s'/s, a) \gamma U(s') \\
 &= R(s) + \max_{a \in A(s)} \sum_{s'} P(s'/s, a) \gamma U(s')
 \end{aligned}$$

$$* \pi(s) = \text{Same, but use } \boxed{\arg \max_{a \in A(s)}} \text{ instead of } \boxed{\max_{a \in A(s)}}$$

② Geometric Sum

$$\rightarrow \sum_{i=1}^n \gamma^i r = \frac{r(1-\gamma^{n+1})}{1-\gamma}, \quad \sum_{i=0}^{\infty} \gamma^i r \leq \frac{R_{\max}}{1-\gamma}$$

* Value Iteration

* حدد القانون حسب شكل الـ R بـ ثابت
 * حدد قيمة Initial و U و خليه بـ 0
 لو صا كفى حاجه
 * اكمل عدد Iterations معين أو لحد فرق أقل من Threshold

* Policy Iteration

* حدد القانون حسب شكل الـ R بـ ثابت
 * حدد Initial Actions سواد من قبل أو Given
 * أ حسب الـ (U) بس المرادي مفيش Max
 فترجع على الـ action الى كملت

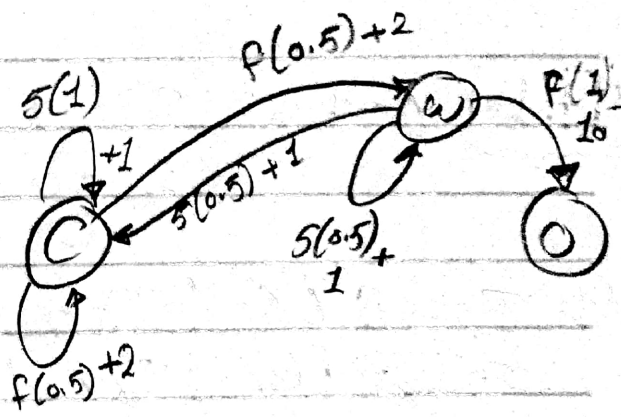
* لو طالع فترك صي معادلات بي بمافيل، حلهم مع بعضي شان تطالع الـ U
 * استعمل معادلات الـ $\pi(s)$ شان تشوف الـ $\arg \max$ وده الـ action الجيد
 * اكمل كده لحد مات Converge

Sh 7

17.4 $\gamma = 0.5$

$$U(s) = \max_{a \in A(s)} R(s, a) + \gamma \sum P(s'/s, a) U(s')$$

* Value Iteration:



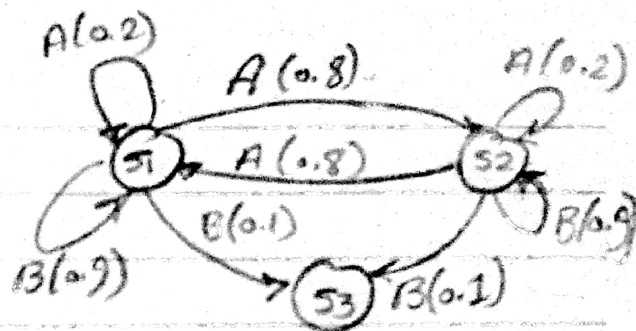
P	$U(C)$	$U(W)$	$U(O)$
0	0	0	0
1	$\text{Max} \begin{bmatrix} 1 + P(C C, 1) \times \gamma U(C) \\ 2 + \begin{bmatrix} P(C C, F) \times \gamma U(C) \\ P(W C, F) \times \gamma U(W) \end{bmatrix} \end{bmatrix}$ $= 2$	$\text{Max} \begin{bmatrix} 1 + \begin{bmatrix} P(W W, 1) \times \gamma U(W) \\ P(C W, 1) \times \gamma U(C) \end{bmatrix} \\ -1 + [P(O W, F) \times \gamma U(O)] \end{bmatrix}$ $= 1$	No Actions $= 0$
2	$\text{Max} \begin{bmatrix} 1 + 1 \times \gamma U(C) \\ 2 + \begin{bmatrix} 0.5 \times \gamma U(C) \\ 0.5 \times \gamma U(W) \end{bmatrix} \end{bmatrix}$ $= 2.75$	$\text{Max} \begin{bmatrix} 1 + \begin{bmatrix} 0.5 \times \gamma U(W) \\ 0.5 \times \gamma U(C) \end{bmatrix} \\ -1 + [1 \times \gamma U(O)] \end{bmatrix}$ $= 1.75$	No Actions $= 0$

* Policy Extraction [No iterations]

$$\pi(s) = \arg \max_{a \in A(s)} R(s, a) + \gamma \sum P(s'/s, a) U(s')$$

$\rightarrow \pi(C) = \arg \max \begin{bmatrix} 1 + 1 \times \gamma U(C) \\ 2 + \begin{bmatrix} 0.5 \times \gamma U(C) \\ 0.5 \times \gamma U(W) \end{bmatrix} \end{bmatrix} = F$
 $\rightarrow \pi(W) = \arg \max \begin{bmatrix} 1 + \begin{bmatrix} 0.5 \times \gamma U(W) \\ 0.5 \times \gamma U(C) \end{bmatrix} \\ -1 + [1 \times \gamma U(O)] \end{bmatrix} = S$
 $\pi(O) = -$

17.1. $\gamma=1, R(s_1)=-1, R(s_2)=-2$
 $R(s_3)=0 \leftarrow \text{Terminal}$



- ① $s_1 \rightarrow B$ [Go to s_3 or stay in s_1], [A] is bad b/c more penalty
 $s_2 \rightarrow A$ [Go to s_1 and stuck there is better], [B] why bad?
 Bec it has lower prob. to succeed!

② Value Iteration $[U(s) = R(s) + \max_{a \in A(s)} [P(s'|s, a) \gamma U(s')]]$

i	$U(s_1)$	$U(s_2)$	$U(s_3)$
0	0	0	0
1	$-1 + \max \begin{bmatrix} 0.2 \times \gamma U(s_1) + \\ 0.8 \times \gamma U(s_2) \\ 0.9 \times \gamma U(s_3) + \\ 0.1 \times \gamma U(s_3) \end{bmatrix}$ $= -1$	$-2 + \max \begin{bmatrix} 0.2 \times \gamma U(s_2) + \\ 0.8 \times \gamma U(s_1) \\ 0.9 \times \gamma U(s_3) + \\ 0.1 \times \gamma U(s_3) \end{bmatrix}$ $= -2$	No Actions $= 0$
2	$-1 + \max \begin{bmatrix} 0.2 \times -1 + 0.8 \times -2 \\ 0.9 \times -1 + 0.1 \times 0 \end{bmatrix}$ $= -1.9$	$-2 + \max \begin{bmatrix} 0.2 \times -2 + 0.8 \times -1 \\ 0.9 \times -2 + 0.1 \times 0 \end{bmatrix}$ $= -3.2$	$= 0$
3	$= -2.71$	$= -4.16$	$= 0$

* Policy Extraction

$$\begin{aligned} \pi(s_1) &= \arg \max \begin{bmatrix} 0.2 \times U(s_1) + 0.8 \times U(s_2) \\ 0.9 \times U(s_1) + 0.1 \times U(s_3) \end{bmatrix} = B \\ \pi(s_2) &= \arg \max \begin{bmatrix} 0.2 \times U(s_2) + 0.8 \times U(s_1) \\ 0.9 \times U(s_2) + 0.1 \times U(s_3) \end{bmatrix} = A \end{aligned}$$

③ Policy Iteration

		S_1	S_2	S_3
0	U	-	-	-
	π	B	B	-
1	U	$-1 + (0.1 \times \gamma U(S_3) + 0.9 \times \gamma U(S_1))$ $= -1$	$-2 + (0.1 \times \gamma U(S_3) + 0.9 \times \gamma U(S_2))$ $= -2$	-
	π	$\arg\max \begin{bmatrix} 0.8U(S_2) + 0.2U(S_1), \\ 0.1U(S_3) + 0.9U(S_1) \end{bmatrix}$ $= B$	$\arg\max \begin{bmatrix} 0.8U(S_1) + 0.2U(S_2), \\ 0.1U(S_3) + 0.9U(S_2) \end{bmatrix}$ $= A$	-
2	U	$-1 + (0.1U(S_3) + 0.9U(S_1))$ $= -1$	$-2 + (0.8U(S_1) + 0.2U(S_2))$ $= -12.5$	-
	π	$\arg\max \begin{bmatrix} 0.8U(S_2) + 0.2U(S_1), \\ 0.1U(S_3) + 0.9U(S_1) \end{bmatrix}$ $= B$	$\arg\max \begin{bmatrix} 0.8U(S_1) + 0.2U(S_2), \\ 0.1U(S_3) + 0.9U(S_2) \end{bmatrix}$ $= A$	-
		Converged	Converged	

④ What if the initial action was 'A' for both?

$$\begin{cases} U(S_1) = -1 + 0.8(\gamma U(S_2)) + 0.2(\gamma U(S_1)) \\ U(S_2) = -2 + 0.8(\gamma U(S_1)) + 0.2(\gamma U(S_2)) \end{cases} \quad \begin{matrix} U(S_1) \rightarrow x \\ U(S_2) \rightarrow y \end{matrix}$$

$$\begin{cases} -0.8x + 0.8y = 1 \\ 0.8x - 0.8y = 2 \end{cases} \quad \begin{matrix} \text{Two // lines} \\ \text{No sol!} \end{matrix}$$

$$\begin{cases} (0.2\gamma - 1)x + 0.8\gamma y = 1 \\ 0.8\gamma x + (0.2\gamma - 1)y = 2 \end{cases} \rightarrow \begin{pmatrix} 0.2\gamma - 1 & 0.8 \\ 0.8 & 0.2\gamma - 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$\text{We have a sol iff } \Delta \neq 0 \rightarrow \frac{(0.2\gamma - 1)^2}{(0.8)^2} \neq 0 \rightarrow \gamma \neq 1 \quad \begin{matrix} \gamma \neq 1 \rightarrow \text{But, } 0 < \gamma < 1 \\ \text{So, No sol!} \end{matrix}$$