

سوال اول :

بله ممکن است

داریم :

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

$\max_x \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')] \leftarrow$  . انجام دهیم . Action  $a$  در مرحله  $i+1$  و

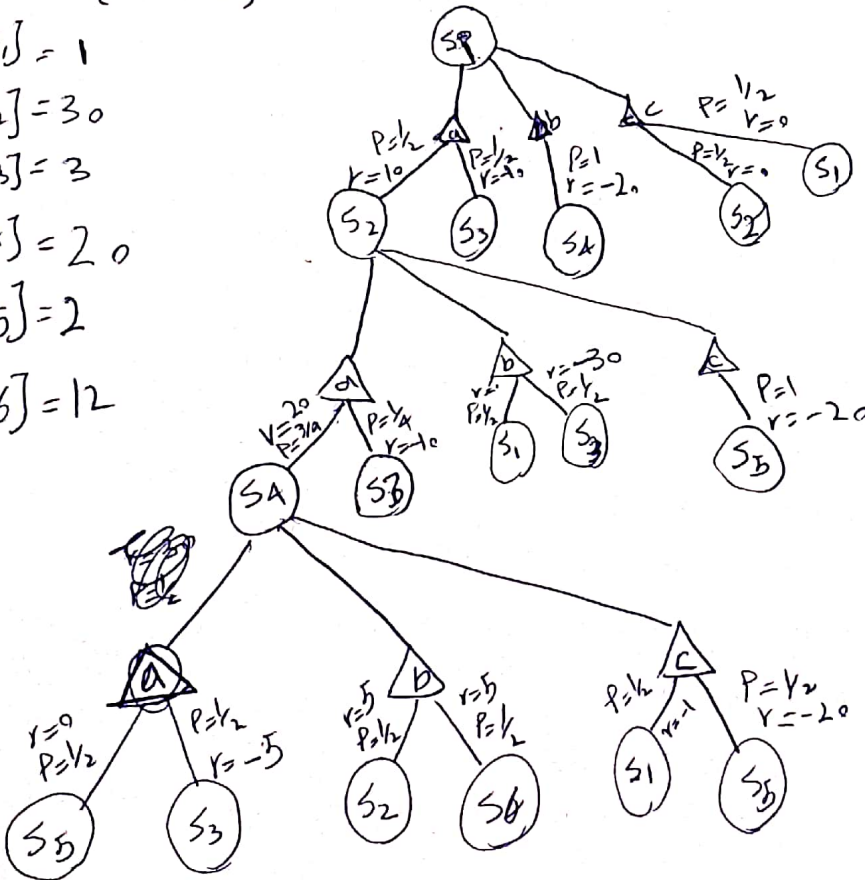
~~max~~

$\max_x \leftarrow \sum_{s'_2} T(s_2, a, s'_2) [R(s_2, a, s'_2) + \gamma V^{\pi_i}(s'_2)]$  . انجام دهیم : Action  $a$  در مرحله  $i+2$  و

اگر فرضی وجود نداشته در مرحله  $i+2$  Action  $a$  انتخاب می‌باشد .  
دلیل اول  $\leftarrow$  میزان Reward  $V^{\pi_i}(s'_2)$  برای  $i+1$  ،  $i+2$  متفاوت است .

Action {a, b, c}

- $V[s_1] = 1$
- $V[s_2] = 30$
- $V[s_3] = 3$
- $V[s_4] = 20$
- $V[s_5] = 2$
- $V[s_6] = 12$



مثال :  $\gamma = 0.9$

در انتخاب بهترین  
بالکس  $\frac{a}{s}$   $\leftarrow$   $s_2$  یا  $s_3$  ؟  
فرض کنید .

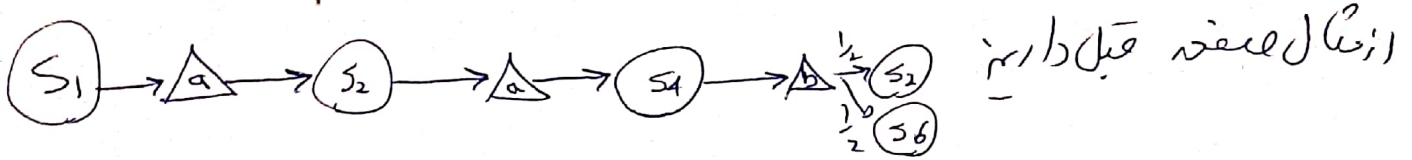
$\rightarrow$  باز هم  $a$  انتخاب  
و فرض کنید  $s_4$  رفتیم .

$\rightarrow$  در اینجا اکس بهترین برای

$\leftarrow$   $\pi$  است .  
بنابراین Policy تعیین کرده است .

علت ۵

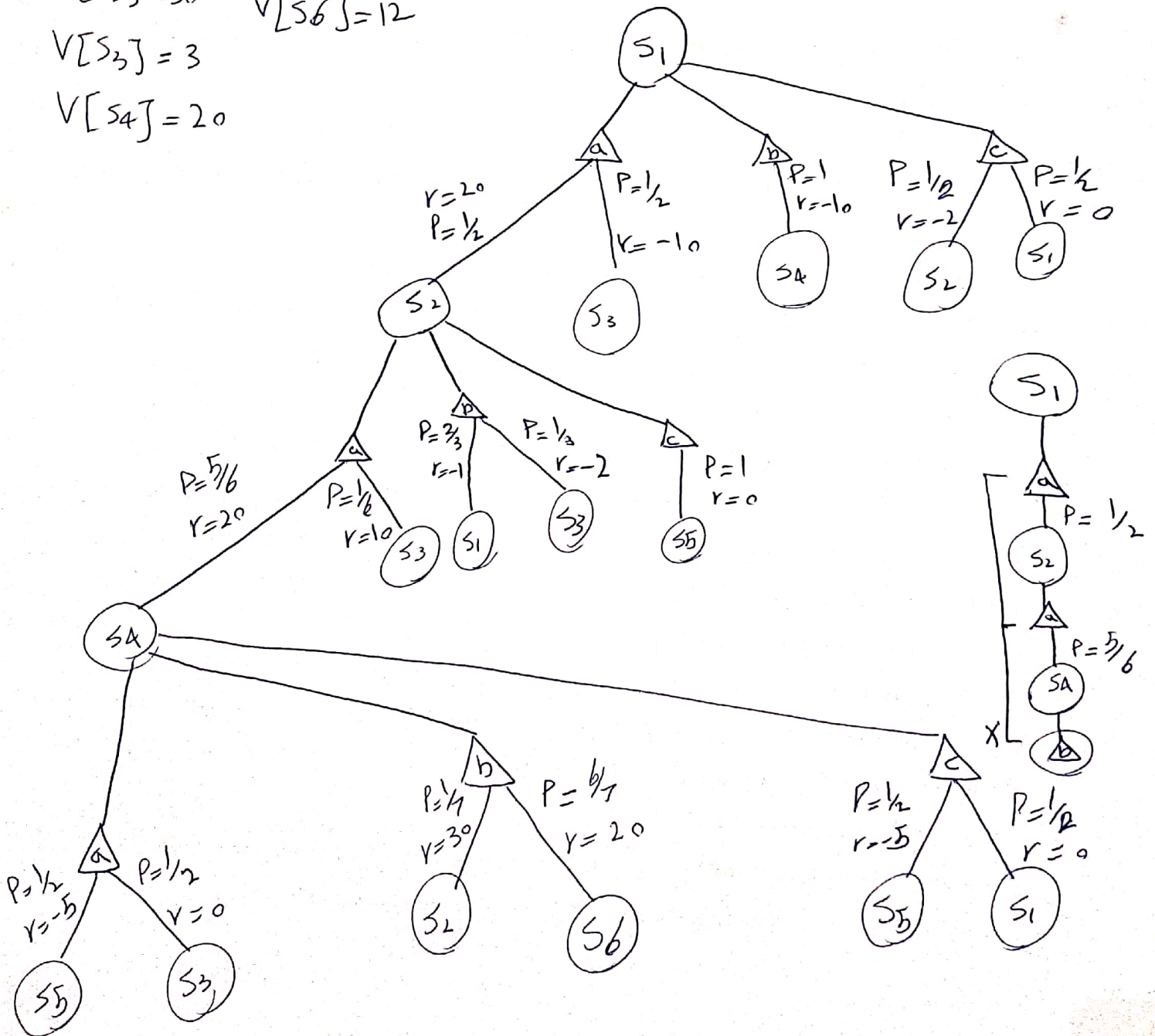
از آنجاده با انتخاب هر اکشن با یک احتمال به یک state می‌رسد.  
 فضای است ممکن است جدید باشد  $\Leftarrow$  transition prob  
 Policy  $\Leftarrow$  می‌تواند تغییر کند



$$\begin{aligned} V[S_1] &= 1 & V[S_5] &= 2 \\ V[S_2] &= 30 & V[S_6] &= 12 \\ V[S_3] &= 3 \\ V[S_4] &= 20 \end{aligned}$$

Actions = {a, b, c}

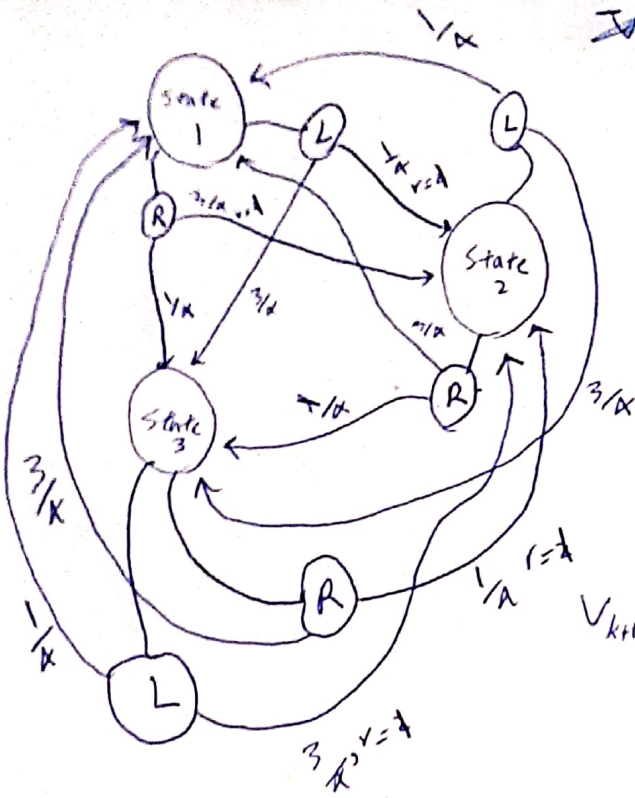
در DMP زیر



به طوری انتخاب می‌شود Policy در هر مرحله به expect و آلتی دارد expect می‌باشد  
 تا به ایاز Transition و مقدار  $V^*$  state است. بنابراین از آنجا که  
 مقدار expect در Policy یکسان در مرحله  $t+1$  و مقدار  $V^*$  است.  
 می‌توان گفت که در مرحله  $t+2$  نیز همان آلتی مقدار expect ماکسیم را  
 دارد است. [حالت مثال اما که زدیم] بنابراین حکم Policy  
 تغییر نکند.



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



$V_k$	$k=0$	1	2
State 1	0	$3/4$	$27/32$
State 2	0	0	$12/32$
State 3	0	$3/4$	$27/32$

$$V_{k+1}(s) = \max_a \left( \sum_a T(s, s', a) (R(s, s') + \gamma V_k(s')) \right)$$

دست کشیدن که منفی است و در هر دو حالت  $V_{s_{i+1}} \text{ reward}(s_i \rightarrow s_2) = 4$

$$V_2[s_1] = \max \left( \frac{1}{4} \times (1 + \frac{1}{2} \times 0) + \frac{3}{4} \times (0 + \frac{1}{2} \times \frac{3}{4}), \frac{3}{4} \times (1 + 0) + \frac{1}{4} \times (\frac{1}{2} \times \frac{3}{4}) \right) = \frac{3 + 24}{32} = \frac{27}{32}$$

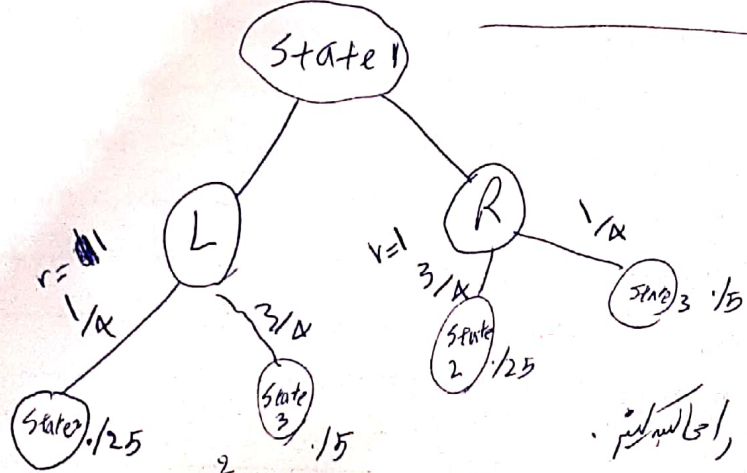
$\frac{1}{4} + \frac{9}{32} = \frac{17}{32}$        $\frac{3}{4} + \frac{3}{32} = \frac{27}{32}$

$$V_2[s_2] = \max \left( \frac{1}{4} \times (\frac{3}{4} \times \frac{1}{2}) + \frac{3}{4} \times \frac{1}{2} \times \frac{3}{4}, \frac{3}{4} \times (\frac{1}{2} \times \frac{3}{4}) + \frac{1}{4} \times \frac{1}{2} \times \frac{3}{4} \right) = \frac{12}{32}$$

$\frac{9}{32} + \frac{9}{32} = \frac{18}{32}$        $\frac{9}{32} + \frac{3}{32} = \frac{12}{32}$

$$V_2[s_3] = \max \left( \frac{1}{4} \times \frac{1}{2} \times \frac{3}{4} + \frac{3}{4} \times (1 + \frac{1}{2} \times 0), \frac{1}{4} \times (1 + \frac{1}{2} \times 1) + \frac{3}{4} \times \frac{1}{2} \times \frac{3}{4} \right) = \frac{27}{32}$$

$\frac{9}{32} + \frac{9}{32} = \frac{18}{32}$        $\frac{1}{4} + \frac{9}{32} = \frac{17}{32}$



باید  $E(L)$  و  $E(R)$  را محاسبه کنیم.

$$E(L) = \frac{3}{4} \times (\frac{1}{2} \times \frac{1}{2}) + \frac{1}{4} \times (1 + \frac{1}{2} \times \frac{3}{4}) = \frac{3}{16} + \frac{9}{32} = \frac{15}{32}$$

$$E(R) = \frac{1}{4} \times (\frac{1}{2} \times \frac{1}{2}) + \frac{3}{4} \times (1 + \frac{1}{2} \times \frac{3}{4}) = \frac{1}{16} + \frac{24}{32} = \frac{25}{32}$$

$(1 + \frac{1}{2} \times \frac{9}{8})$

بنابراین اکشن  $R$  بهینه‌تر است.

با انتخاب  $R$  با احتمال  $\frac{3}{4}$  به state 2 و  $\frac{1}{4}$  به state 1 = Reward 1 داریم. یا به احتمال  $\frac{1}{4}$  action  
 به state 3 داریم. بنابراین با توجه به معیار expectimax و policy این انتخاب  $R$  است.



$V \backslash k$	0	1	2
state1	0	$\frac{3}{4}A$	$\frac{27}{32}$
state2	0	0	$\frac{3}{8}$
state3	0	$\frac{3}{4}A$	$\frac{27}{32}$

خلاصه  
 ان ←

ب ← Action  $R$  زیرا معیار expectimax آن بیشتر است.