

$$Q(s,a) = \sum_{s'} P_{a|s}(s'|s) [r(s,a,s') + \gamma V(s')]$$

$$V(s,a) = \max_a Q(s,a)$$

Q & A

①  $\gamma = 1.0$

$$Q(\text{Messi}, \text{pass}) = P_{\text{pass}}(\text{Suarez} | \text{Messi}) [r(\text{Messi}, \text{pass}, \text{Suarez}) + \gamma V(\text{Suarez})]$$

$$= 1 \times [-1 + 1 \times (-1.2)] = -2.2$$

$$Q(\text{Messi}, \text{shot}) = P_{\text{shot}}(\text{Goal} | \text{Messi}) [r(\text{Messi}, \text{shot}, \text{Goal}) + \gamma V(\text{Goal})]$$

$$+ P_{\text{shot}}(\text{Suarez} | \text{Messi}) [r(\text{Messi}, \text{shot}, \text{Suarez}) + \gamma V(\text{Suarez})]$$

$$= 0.8 \times [-2 + 1 \times (-1.2)] + 0.2 \times (-2 + 1 \times 1) = -2.76$$

$$= -2.76$$

$$= -2.76$$

$$V(\text{Messi}, a) = \max(Q(\text{Messi}, \text{pass}), Q(\text{Messi}, \text{shot}))$$

$$= \max(-2.2, -2.76) = -2.2$$

$\therefore \text{pass}$

$$\textcircled{2} Q(\text{Su}, \text{pass}) = P_{\text{pass}}(\text{Messi} | \text{Su}) [r(\text{Su}, \text{pass}, \text{Messi}) + \gamma V(\text{Messi})] = -3$$

$$Q(\text{Su}, \text{shot}) = P_{\text{shot}}(\text{Messi} | \text{Su}) [r(\text{Su}, \text{shot}, \text{Messi}) + \gamma V(\text{Messi})]$$

$$+ P_{\text{shot}}(\text{Goal} | \text{Su}) [r(\text{Su}, \text{shot}, \text{Goal}) + \gamma V(\text{Goal})]$$

$$= 0.4 \times [-2 + 1 \times (-2)] + 0.6 \times [-2 + 1 \times 1] = -2.2$$

$$V(\text{Goal}) = Q(\text{Goal}, \text{return}) = P_{\text{return}}(\text{Messi} | \text{Goal}) [r(\text{Goal}, \text{return}, \text{Messi}) + \gamma V(\text{Messi})]$$

$$= 1 \times [2 + 1 \times (-2)] = 0$$

③  $\gamma = 0.8$  policy -

$$V^{\pi}(\text{Messi}) = Q^{\pi}(\text{Messi}, \text{pass}) = P_{\text{pass}}(\text{Su} | \text{Messi}) [r(\text{Messi}, \text{pass}, \text{Su}) + \gamma V^{\pi}(\text{Su})] = 1 \times [-1 + 0.8 V^{\pi}(\text{Su})] = 0.8 V^{\pi}(\text{Su}) - 1$$

$$V^{\pi}(\text{Su}) = Q^{\pi}(\text{Su}, \text{pass}) = 1 \times [-1 + 0.8 V^{\pi}(\text{Messi})] = 0.8 V^{\pi}(\text{Messi}) - 1$$

$$V^{\pi}(\text{Goal}) = Q^{\pi}(\text{Goal}, \text{return}) = 1 \times [2 + 0.8 V^{\pi}(\text{Messi})] = 2 + 0.8 V^{\pi}(\text{Messi})$$

$$\begin{cases} a = \gamma b - 1 \\ b = \gamma a - 1 \\ c = 2 + \gamma a \end{cases}$$

$$a = \gamma(\gamma a - 1) - 1 = \gamma^2 a - \gamma - 1 = a$$

$$(\gamma^2 - 1)a = \gamma + 1$$

$$a = \frac{\gamma + 1}{(\gamma^2 - 1)} = \frac{1}{\gamma - 1}$$

$$\therefore \begin{cases} a = \frac{1}{\gamma - 1} = -5 \\ b = \frac{1}{\gamma - 1} = -5 \end{cases}$$

$$b = \frac{1}{\gamma - 1} = -5$$

$$c = 2 + 0.8 \times (-5) = -2$$

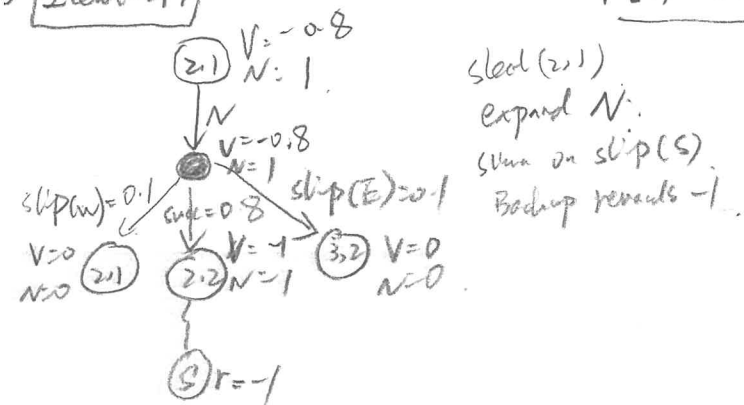
$$\therefore 0.8 \times (-6) + 0.2 \times (-36) = -4.8 + (-0.72) = -5.52$$

$$= 0.4 \times [-2 + 0.8 \sqrt{11}] + 0.6 \times [-2 + 0.8 \sqrt{11}]$$

WS & MCS

$$V_L(s) = \max_a \sum P_a(s'|s) [r + \gamma V(s')]$$

Iteration 1



sled (2,1)  
 expand N.  
 turn on slip(s).  
 Backup reveals -1.

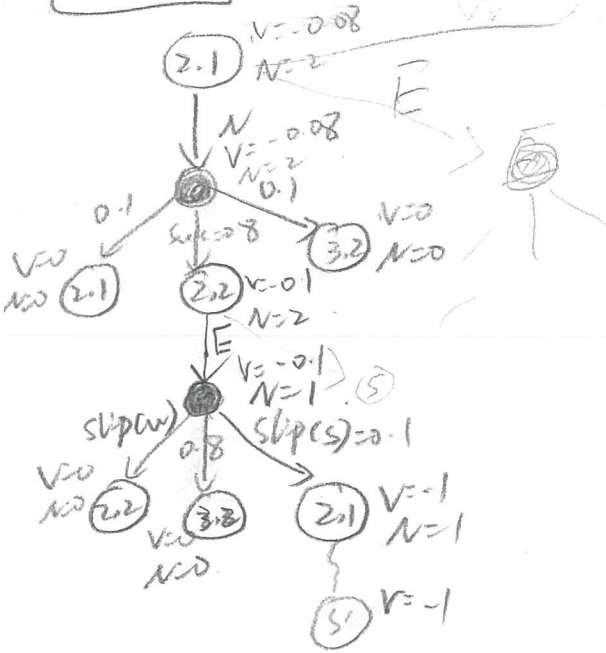
$$Q(2,1, E) = -0.8$$

②  ~~$w = 0.8$~~   $Q(2,1), w = 0.8$ .

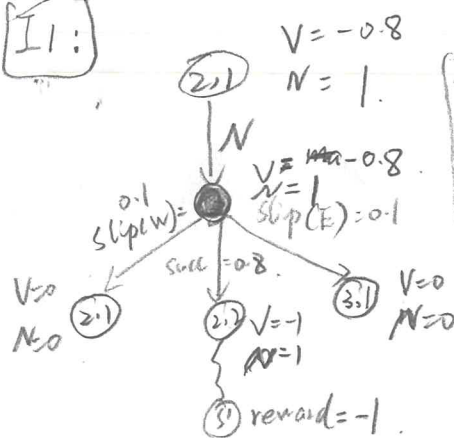
$$UCT = \operatorname{argmax}_S Q(s, a) + 2 \cdot c_p \cdot \sqrt{2 \ln N(s)}$$

$$\pi(s) = \begin{cases} E: 0.8 + \sqrt{\frac{2 \ln(5)}{1}} \\ N: 0.08 + \sqrt{\frac{2 \ln(5)}{3}} \\ W: 0.8 + \sqrt{\frac{2 \ln(5)}{1}} \\ S: \infty \end{cases}$$

② Zterolb n z



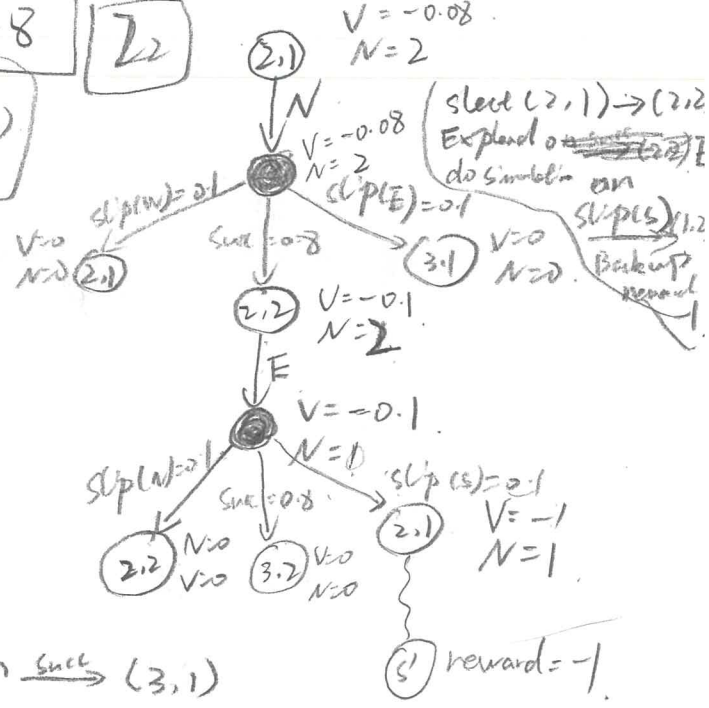
I1:



$$V(s) = 0.1 \times (0 + 0.1 \times 0) + 0 + 0.8 \times (-1 + 0.1 \times 0) = -0.8$$

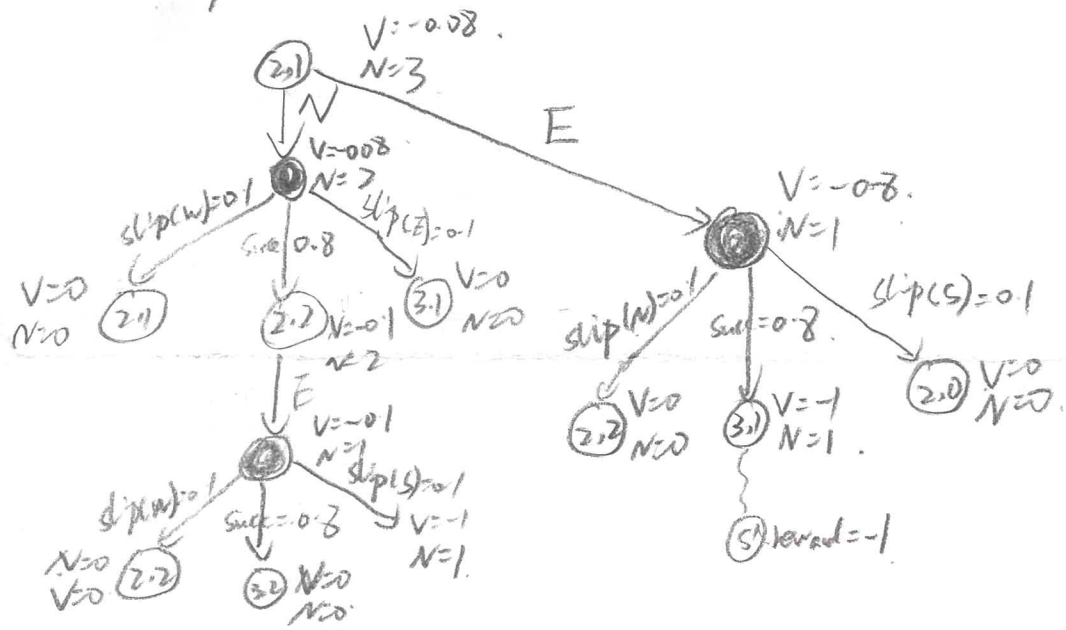
Week 8

I2



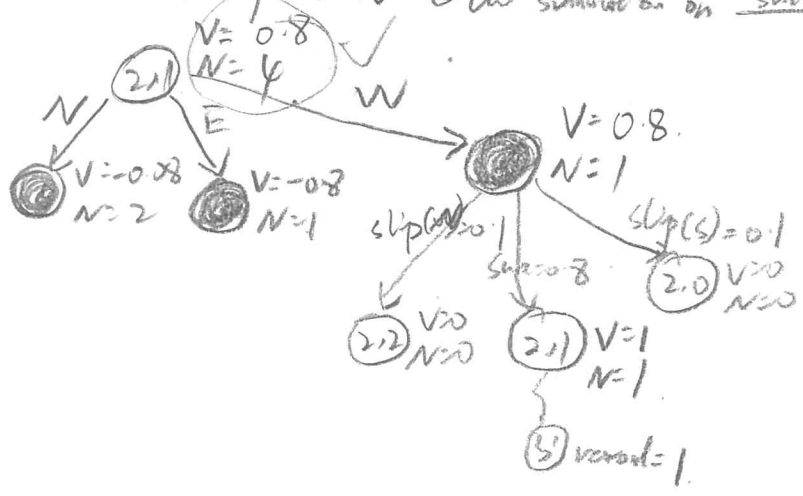
I3

- ① select (2,1) ② expand E ③ do simulation on succ -> (3,1) ④ Backup the rewards -1.



I4

- ① Select (2,1) ② expand W ③ do simulation on succ -> (2,1) ④ Backup the reward -1.



④ But ~~un~~ it reads = 1



$$Q(12, 1, \bar{E}) = -0.8$$

$$Q(z, s) = 0$$



$$\alpha = 0.4 \quad \gamma = 0.9$$

Week 9

$$\begin{aligned} [2] \quad Q(S, \text{pass}) &= Q(S, \text{pass}) + \alpha [r(S, \text{pass}, M) + \gamma \max_{a' \in A(M)} Q(M, a') - Q(S, \text{pass})] \\ &= -0.7 + 0.4 [-1 + 0.9 \times (-0.4) - (-0.7)] \\ &= -0.7 + 0.4 (-1 - 0.36 + 0.7) = -0.7 - 0.264 = \underline{-0.964} \end{aligned}$$

$$\begin{array}{r} -0.66 \\ \times 0.4 \\ \hline 0.264 \end{array}$$

[3] SARSA:

$$\begin{aligned} Q(S, P) &= Q(S, P) + \alpha [r(S, \text{pass}, M) + \gamma Q(M, S) - Q(S, P)] \\ &= -0.7 + 0.4 [-1 + 0.9 \times (-0.8) + 0.7] \\ &= -0.7 + 0.4 (-1 - 0.72 + 0.7) = \underline{-1.108} \end{aligned}$$

$$\begin{array}{r} -1.02 \\ \times 0.4 \\ \hline -0.408 \end{array}$$

[4]  $\star$

$$Q(S, a) = Q(S, a) + \alpha [G_t^n - Q(S, a)] \quad G_t = \sum_{i=t+1}^T \gamma^i r_i$$

$$Q(S, a) + \alpha [G_t + \gamma^n Q(S, a') - Q(S, a)]$$

$$Q(S_{\text{Suarez}}, \text{Pass}) = Q(S_{\text{Suarez}}, \text{Pass}) + \alpha [G_{\text{Suarez}}^3 + \gamma Q(M, P) - Q(S, P)]$$

$$= -0.7 + 0.4 [G_{\text{Suarez}}^3 + 0.9 \times (-0.4) + 0.7]$$

$$= -0.7 + 0.4 [-1.18 + (-0.36) + 0.7] = -0.7 + 0.4 \times (-0.7716) = \underline{-1.00864}$$

$$G_{\text{Suarez}}^3 = r(S, P) + \gamma r(M, S) + \gamma^2 r(S_{\text{Suarez}}, R)$$

$$= -1 + 0.9 \times (-2) + 0.9^2 \times 2 = -1 - 1.8 + 1.62 = \underline{-1.18}$$

action (S, P)  
 ① (S, P)  
 ② (M, S)  
 ③ (A, P)  
 ④ (S, R)  
 (Scored, R)

$$\begin{array}{r} 0.81 \\ \times 0.9 \\ \hline 0.729 \end{array}$$

$$\begin{array}{r} -0.7716 \\ \times 0.4 \\ \hline -0.30864 \\ + 0.7 \\ \hline 1.00814 \end{array}$$

$$\begin{array}{r} -0.18 \\ \times 0.81 \\ \hline -0.324 \\ + 0.7 \\ \hline -0.324 \\ + 0.376 \\ \hline 1.18 \\ + 1.18 \\ \hline 2.36 \\ + 0.324 \\ \hline 1.504 \\ + 1.18 \\ \hline 2.684 \\ + 0.2916 \\ \hline 1.4716 \end{array}$$

$$\begin{array}{r} 0.729 \\ \times 0.4 \\ \hline 0.2916 \end{array}$$

# [Week 9]

[2]  $Q(s,a) = Q(s,a) + \alpha [r(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$

$Q(S_u, \text{pass}) = -0.7 + 0.4 [-1 + 0.9 \times \text{max } Q(M, \text{pass}) + 0.7]$   
 $= -0.7 + 0.4 [-1 - 0.36 + 0.7] = -0.7 - 0.264 = -0.964$

[3]  $Q(S, \text{pass}) = -0.7 + 0.4 [-1 + 0.9 \times (-0.8) + 0.7]$   
 $= -0.7 + 0.4 [-1 - 0.72 + 0.7] = -0.7 - 0.408 = -1.108$

[4]  $Q(S_u, \text{pass})$   $r_1$   
 $Q(M, \text{shoot})$   $r_2$   
 $Q(S_{\text{coord}}, \text{relve})$   $r_3$   
 $Q(M, \text{pass})$   $r_4$

3-sleep SARSA

$Q(s,a) = Q(s,a) + \alpha [G_t + \gamma^3 Q(s',a') - Q(s,a)]$

$G_t = r_1 + \gamma r_2 + \gamma^2 r_3 = -1 + 0.9 \times (-2) + 0.9^2 \times 2 = -1 - 1.8 + 1.62 = -1.18$

$Q(S_u, \text{pass}) = -0.7 + 0.4 [-1.18 + 0.9 \times (-0.4) + 0.7]$   
 $= -0.7 + 0.4 [-1.18 - 0.36 + 0.7] = -0.7 - 0.408 = -1.108$