

6.1 TD PREDICTION

NON STATIONARY EVERY-VISIT MONTI CARLO:

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

→  $\alpha$ : STEP SIZE: IF  $\alpha = 1$ , THEN WE SET ESTIMATE EQUAL TO LAST MEASUREMENT

→ THE ABOVE FORMULA IS APPLIED WHEN EPISODE FINISHES

ONE-STEP TD, OR TD(0) CORRECTION

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

NOTE THE TEMPORAL DIFFERENCE BETWEEN THE TWO

BETTER ESTIMATE, USING ONE EXTRA DATA POINT,  $R_{t+1}$  INITIAL ESTIMATE

(WE COULD USE MORE DATA POINTS, THIS IS THE TD(n) METHOD). PSEUDO CODE:

INPUT: POLICY  $\pi$  &  $J$ .

PARAMETER: STEP SIZE  $\alpha \in (0, 1]$

INITIALIZE  $V(S)$   $\forall S \in \mathcal{S}$ , BUT  $V(\text{TERMINAL}) = 0$

LOOP FOR EACH EPISODE:

INITIALIZE  $S$

LOOP FOR EACH STEP OF EPISODE:

$A \leftarrow$  ACTION GIVEN BY  $\pi$  FOR  $S$

TAKE ACTION, OBSERVE  $R, S'$

$$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$$

$S \leftarrow S'$

UNTIL  $S$  TERMINAL

- LET US COMPARE
- 1) DP (DYNAMIC PROGRAMMING)
  - 2) MC (MONTE CARLO)
  - 3) TD (TEMPORAL DIFFERENCES)

$$Q_{\pi}(s) \triangleq E_{\pi} [ G_t | S_t = s ] \leftarrow \text{MONTE CARLO ESTIMATES THIS}$$

$$= E_{\pi} [ R_{t+1} + \gamma G_{t+1} | S_t = s ]$$

$$= E_{\pi} [ R_{t+1} + \gamma Q_{\pi}(S_{t+1}) | S_t = s ]$$

TD ESTIMATES BOTH

DP "ESTIMATES" THIS

IMPORTANT QUANTITY:

TD ERROR

TD ERROR:  $\delta_t \triangleq R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$

MONTE CARLO ERROR:

$$G_t - V(S_t)$$

THE TWO ARE CONNECTED. ASSUMING THE ESTIMATED VALUES DO NOT CHANGE,

$$G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1})$$

$$= \delta_t + \gamma (G_{t+1} - V(S_{t+1}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 (G_{t+2} - V(S_{t+2}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2} + \dots + \gamma^{T-t} (G_T - V(S_T))$$

$$= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$

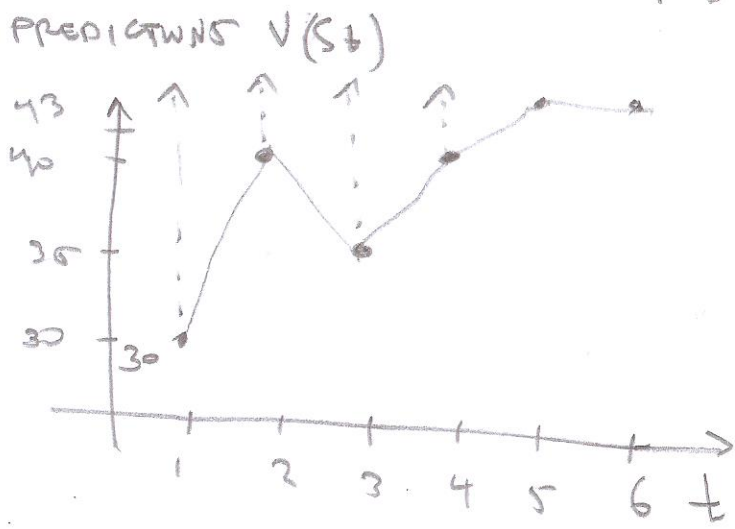
0 0

EXAMPLE 6.1. DRIVING HOME

LET US CONSIDER ONE EPISODE OF DRIVING HOME, ALONG WITH A SPECIAL POLICY (I.E. ROUTE TAKEN, ETC.)

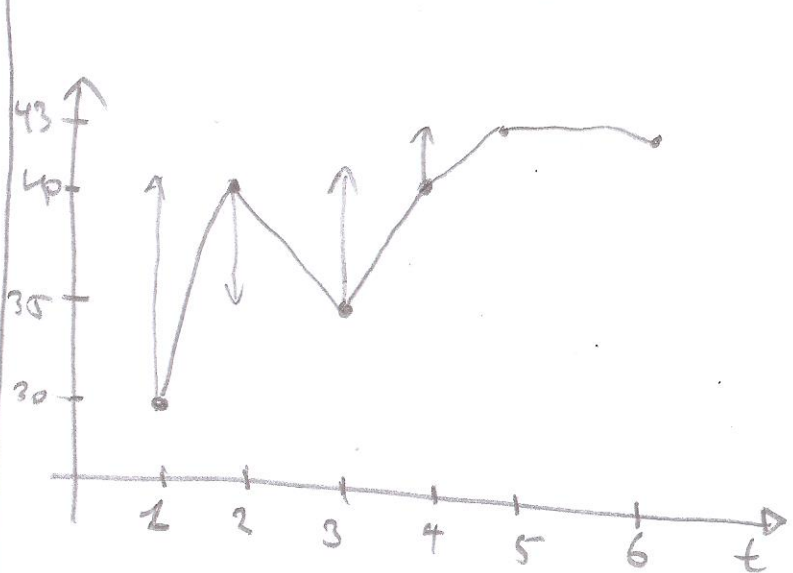
STATE	ELAPSED TIME (REWARDS)	PREDICTED TIME TO GO	PREDICTED TOTAL TIME
1) LEAVING OFFICE FRIDAY AT 6	0	30	30
2) HEAR CAR, RAINING* (* THROUGH ANOTHER STATE WHERE YOU VISIT CAR BUT IT DOES NOT RAIN)	5	30	40
3) EXITING HIGHWAY	20	15	35
4) 2 MAIN ROAD BEHIND TRUCK	30	10	40
5) ENTERING HOME STREET	40	3	43
6) ARRIVING HOME	43	0	43

MONTE CARLO WITH  $\alpha=1$



AT THE END, THESE ALL BECOME EQUAL TO 43

TD(0) WITH  $\alpha=1$

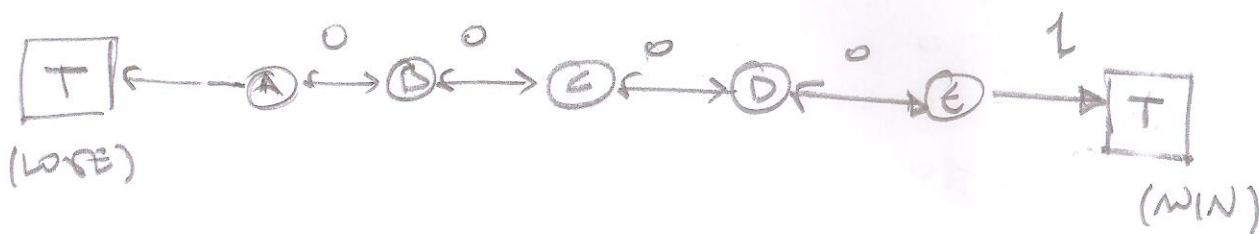


AFTER EACH ITERATION, THESE BECOME EQUAL TO SUBSEQUENT PREDICTION

## 6.2 ADVANTAGES OF TD PREDICTION METHODS

- 1) W.R.T DP, NO NEED FOR MODEL
- 2) W.R.T. MC, NO NEED TO WAIT FOR EPISODE TO END
- 3) THEY CONVERGE (IE.  $V_t(s) \rightarrow U_{\pi}(s)$ )
  - 1) IN THE MEAN, IF  $\alpha$  IS FIXED AND SUFFICIENTLY SMALL
  - 2) WITH PROBABILITY 1, IF  $\alpha \rightarrow 0$  IN USUAL MANNER (2.7)
- 4) TD IS USUALLY FASTER.

### EXAMPLE 6.2 RANDOM WALK (GAMBLER'S ROUIN)



WHAT ARE PROBABILITIES  $P_A, P_B, P_C, P_D, P_E$  STARTING AT RESPECTIVE STATE WE WIN? THAT OBSERVE THAT

$$P_A = \frac{1}{2} \cdot 0 + \frac{1}{2} P_B$$

$$P_B = \frac{1}{2} P_A + \frac{1}{2} P_C$$

$$P_C = \frac{1}{2} P_B + \frac{1}{2} P_D$$

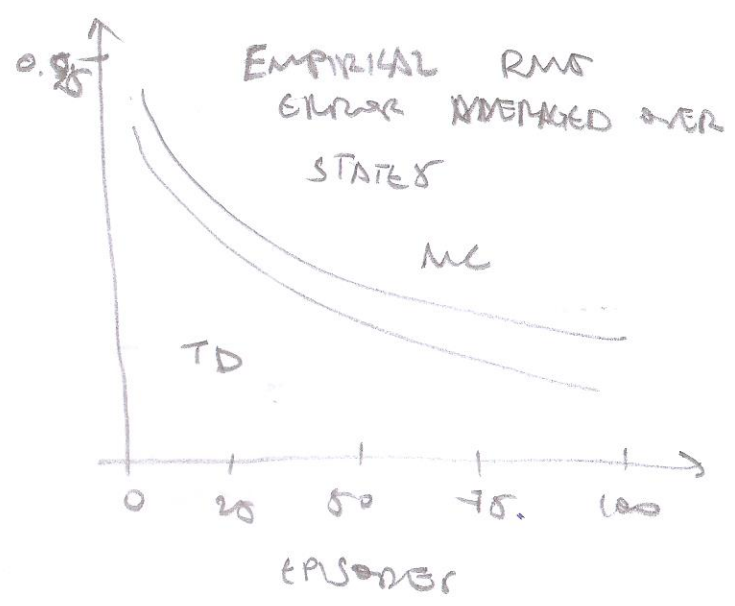
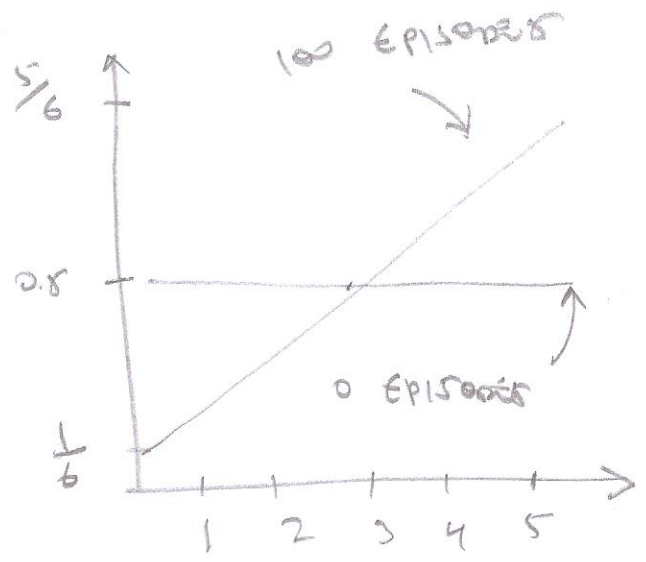
$$P_D = \frac{1}{2} P_E + \frac{1}{2} P_C$$

$$P_E = \frac{1}{2} \cdot 1 + \frac{1}{2} P_D$$

$$\Rightarrow \left. \begin{aligned} P_A &= \frac{1}{6}, & P_B &= \frac{2}{6}, & P_C &= \frac{3}{6} \\ P_D &= \frac{4}{6}, & P_E &= \frac{5}{6} \end{aligned} \right\}$$

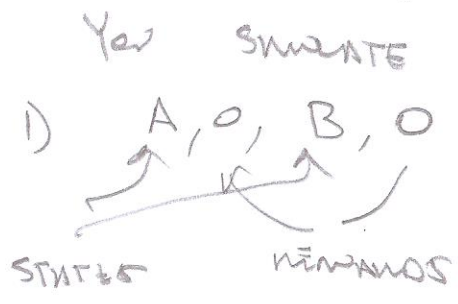
EASY TO GENERALIZE FOR BIASED WALKS AND  $n$  STATES.

EXPLAINS ADVANTAGE OF CASINO



SO, WHY IS IT FASTER? IT USES FURTHER INFO.  
 ALSO, CONSIDER NEXT EXAMPLE:

EXAMPLE 6.4



- 8 EPISODES:
- 1) A, 0, B, 0
  - 2) B, 1, , 3) B, 1, 4) B, 1
  - 5) B, 1 6) B, 1 7) B, 1 8) B, 0

(AFTER B, WE ENTER TERMINAL STATE)

WHAT IS YOUR ESTIMATE OF  $U_{\pi}(A)$  AND  $U_{\pi}(B)$ ?  
 ASSUME  $\gamma = 1$

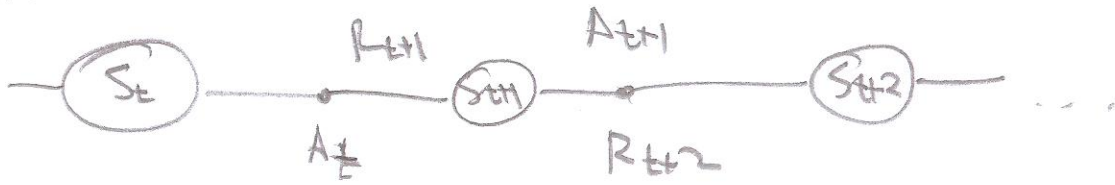
→ IF YOU SAY  $U_{\pi}(B) = 0.75$ ,  $U_{\pi}(A) = 0.75$  THEN  
 YOU SAY WHAT TD(0) ACHIEVES, IF WE TRAIN THE ALGORITHM REPEATEDLY ON ABOVE 8 EPISODES

→ IF YOU SAY  $U_{\pi}(B) = 0.75$ ,  $U_{\pi}(A) = 0$  THEN  
 YOU SAY WHAT MC ACHIEVES, IF WE TRAIN THE ALGORITHM REPEATEDLY ON ABOVE 8 EPISODES

IF MODEL IS MARKOVIAN, FIRST RESULT IS BETTER

# 6.4. SARSA: ON-POLICY TD CONTROL

CONSIDER A SEQUENCE OF STATE-ACTION PAIRS IN AN EPISODE:



THE STATE UPDATE FOR  $V_t(S_t)$  NOW HOLDS FOR  $Q(S_t, A_t)$ :

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

(IF  $S_{t+1}$  = TERMINAL,  $Q(S_{t+1}, A_{t+1}) = 0$ )

UPDATE USES  $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$ , AND  $\alpha$  IS CALLED SARSA

## SARSA (ON-POLICY TD CONTROL) FOR ESTIMATING $Q \approx q^*$

ALGORITHM PARAMETERS: STEP SIZE  $\alpha \in (0, 1]$ , SMALL  $\epsilon > 0$   
 INITIALIZE  $Q(s, a) \forall s \in S^+, a \in A(s)$ , ARBITRARILY BUT  $Q(\text{TERMINAL}, *) = 0$

LOOP FOR EACH EPISODE:

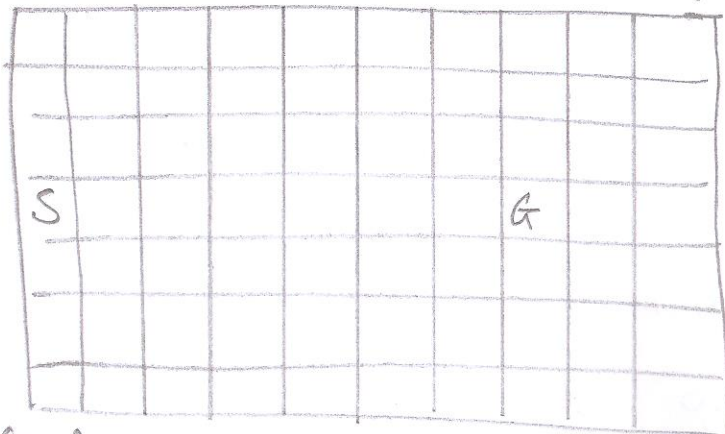
```

INITIALIZE S
CHOOSE A FROM S USING Q (E.G.  $\epsilon$ -GREEDY)
LOOP FOR EACH STEP OF EPISODE:
    TAKE ACTION A, OBSERVE R, S'
    CHOOSE A' FROM S' USING Q (E.G.  $\epsilon$ -GREEDY)
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S', A \rightarrow A'$ 
UNTIL S TERMINAL
    
```

EXAMPLE 6.5

WINDY GRIDWORLD

(9,9)

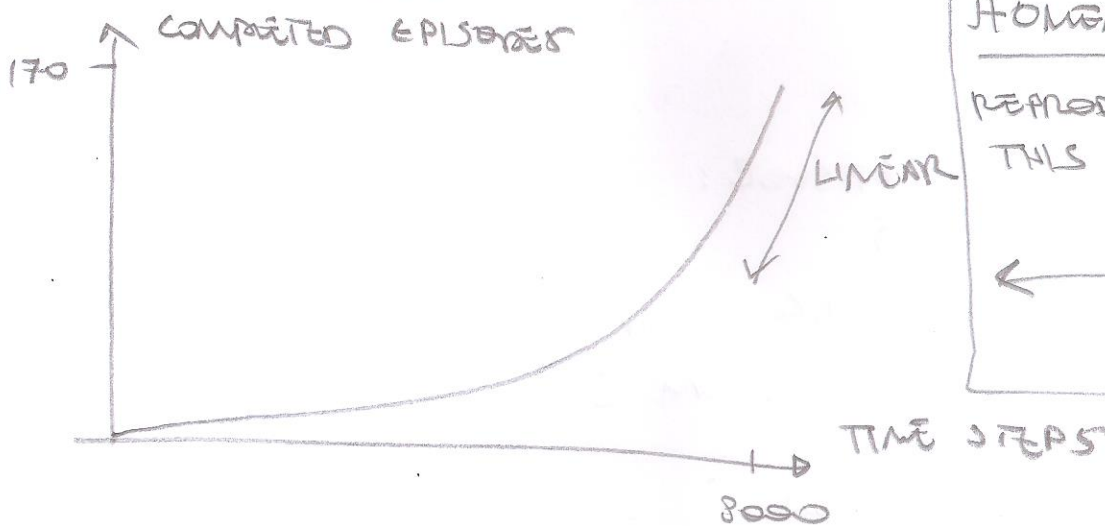


- ALL REWARDS ARE -1
- EPISODIC TASK: WE FINISH AT G
- WE FINISH AT S
- $\epsilon = 0.1$ ,  $\alpha = 0.5$
- INITIAL  $Q(s,a) \forall s,a$

STATE (0,0)

0 0 0 1 1 2 2 1 0

OBSERVE: TRIVIAL TO SOLVE WITH DIJKSTRA'S ALGORITHM, IF WE HAD MODEL.



HOMWORK #6

REPRODUCE THIS FIGURE

6.5 Q-LEARNING: OFF-POLICY TD-CONTROL

CONSIDER ALTERNATIVE UPDATE RULE:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

OBSERVE THAT:

$$(3.20) \quad q^*(s, a) = \sum_{s', r} P(s', r | s, a) \left[ r + \gamma \max_{a'} q^*(s', a') \right]$$

THEFORE, Q-LEARNING Tries TO DISCOVER OPTIMAL POLICY. THEREFORE IT IS AN OFF-POLICY ALGORITHM, BECAUSE THERE ARE TWO POLICIES:

- 1) TARGET POLICY: THE OPTIMAL ONE
- 2) BEHAVIOR POLICY:  $\epsilon$ -GREEDY, & ANY OTHER THAT DOES SOME EXPLORING

Q-LEARNING CAN ALSO BE SHOWN TO CONVERGE.

Q-LEARNING (OFF-POLICY TD CONTROL) FOR ESTIMATING  $J^*$

ALGORITHM PARAMETERS: STEP SIZE  $\alpha \in (0, 1]$ , SMALL  $\epsilon > 0$   
 INITIALIZE  $Q(s, a)$ , FOR ALL  $s \in S^+$ ,  $a \in A(s)$  ARBITRARILY EXCEPT THAT  $Q(\text{TERMINAL}, \cdot) = 0$

LOOP FOR EACH EPISODE:

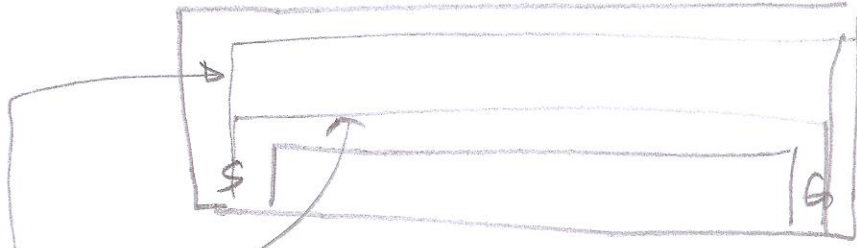
```

INITIALIZE S
LOOP FOR EACH STEP OF EPISODE:
  CHOOSE A FROM S USING Q AND  $\epsilon$ -GREEDY POLICY (OR ANOTHER)
  TAKE ACTION A, OBSERVE R, S'
   $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_{a'} Q(s', a') - Q(s, A)]$ 
   $S \leftarrow S'$ 
UNTIL S TERMINAL
  
```

Q-LEARNING CAN PERFORM WORSE THAN SARSA.

CONSIDER THE CLIFF WALKING EXPERIMENT





Q LEARNING DISCOVERS THE OPTIMAL ROUTE AND USES IT, AND ONE IN A WHILE IT FALLS IN THE CLIFF

BUT SARSA LEARNS TO STAY AWAY FROM THE CLIFF!



6.6 EXPECTED SARSA

WE MODIFY THE UPDATE RULE OF Q-LEARNING TO TAKE AVERAGE, INSTEAD OF MAXIMUM:

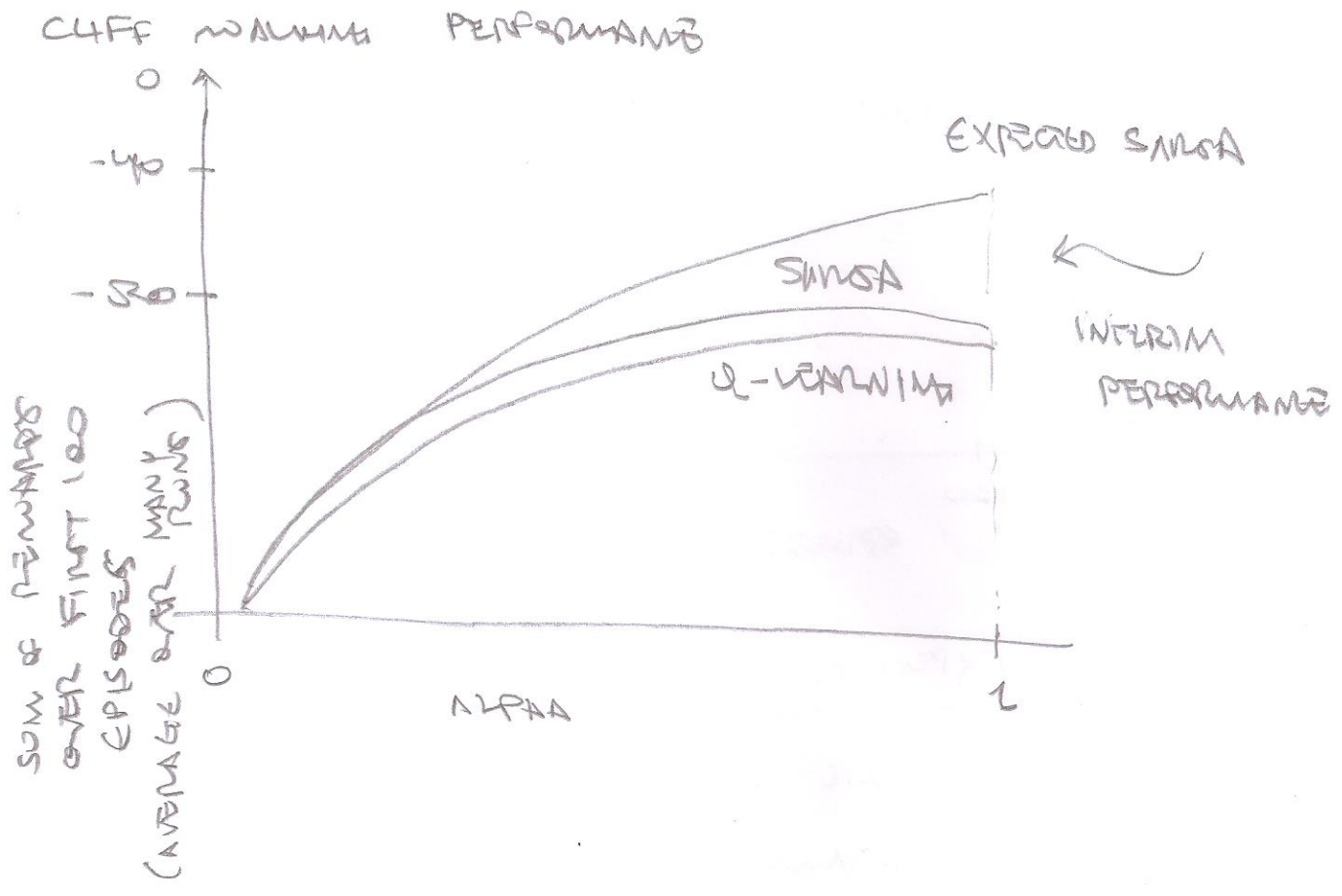
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma E_{\pi} [Q(S_{t+1}, A_{t+1}) | S_{t+1}] - Q(S_t, A_t)]$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_{a'} \pi(a' | S_{t+1}) Q(S_{t+1}, a') - Q(S_t, A_t)]$$

(OBSERVE:  $q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q(s', a')]$ )

SO THIS SUCENT<sup>2</sup> MOVES DETERMINISTICALLY IN THE EXPECTED DIRECTION TOWARDS WHICH SARSA MOVES, AND SO IT CALLED EXPECTED SARSA

BECAUSE IT HAS LESS VARIANCE THAN SARSA, IT GENERALLY PERFORMS BETTER THAN SARSA



**HOMEWORK #7**

REPRODUCE THIS FIGURE, ONLY FOR THE INTERIM CURVES.

BONUS IF YOU CAN EXPLAIN WHY THE ASYMPTOTIC PERFORMANCE FOR SARSA IS BELOW THE INTERIM PERFORMANCE

USE  $\epsilon$ -GREEDY WITH  $\epsilon = 0.1$