

③ Policy Iteration

```

Initialize V(s) arbitrarily;
Randomly initialize policy π(s);
Δ<=0;
while Δ < 0 do
    for each s ∈ S do
        v ← V(s);
        v(s) ←  $\sum_{a \in A} p(a,r|s,a) [r + γV(a)]$ ;
        Δ ← max |Δ, |v - V(s)||;
    end
    Policy-stable ← true;
    for each s ∈ S do
        old-action ← π(s);
        π(s) ← argmax  $\sum_a p(a,r|s,a) [r + γV(a)]$ ;
        if old-action ≠ π(s) then
            Policy-stable ← false;
    end
    end
    if Policy-stable then
        return V as V* and π as π*;
    else
        go to Policy evaluation;
    end
end

```

Policy Evaluation

Policy Improvement

④ Value Iteration

```

Initialize V(s) arbitrarily, except V(terminal);
V(terminal) ← 0;
Δ<=0; (Loop until convergence)
while Δ < 0 do
    for each s ∈ S do
        V ← V(s);
        V(s) ← max  $\sum_{s',a} p(s',r|s,a) [r + γV(s')]$ ;
        Δ ← max |Δ, |V - V(s)||;
    end
    return π, V, π(V);
end

```

⑤ On-policy MC control

```

Initialize, for all s ∈ S, a ∈ A(s);
Q(s,a) ← arbitrary
Returns(s,a) ← empty list
π(a|s) ← an arbitrary ε-soft policy
Repeat forever:
    Generate an episode using π
    For each pair (s,a) appearing in the episode :
        G ← return following the first occurrence of s,a
        append G to Returns(s,a)
        Q(s,a) ← average>Returns(s,a))
    For each S in episode :
        Atk ← argmax  $\sum_a Q(s,a)$ 
    For all a ∈ A(s) :
        π(a|s) ←  $\begin{cases} 1-\epsilon + \frac{\epsilon}{|A(s)|} & \text{if } a = Atk \\ \frac{\epsilon}{|A(s)|} & \text{if } a \neq Atk \end{cases}$ 

```

⑥ off-policy MC prediction

```

Initialize, for all s ∈ S, a ∈ A(s):
Q(s,a) ← arbitrary
C(s,a) ← 0
Repeat forever:
    b ← any policy with coverage of π
    Generate an episode using b:
        So, A0, R1, ..., St-1, At-1, Rt, St
    G ← 0
    W ← 1
    For t=T-1, T-2, ... down to 0 .
        G ← rG + R
        C(st,At) ← C(st,At) + 1
        Q(st,At) ← Q(st,At) +  $\frac{W}{C(st,At)} [G - Q(st,At)]$ 
        W ← W *  $\frac{1}{C(st,At)}$ 
    If W=0 then exit for loop

```

⑦ off-policy MC control

```

Initialize, for all s ∈ S, a ∈ A(s):
Q(s,a) ← arbitrary
C(s,a) ← 0
T(s) ← argmax Q(s,a)
Repeat forever:
    b ← any soft policy
    Generate an episode using b: So, A0, R1, ..., St, At, Rt, St
    G ← 0
    W ← 1
    For t=T-1, T-2, ... down to 0 :
        G ← rG + R
        C(st,At) ← C(st,At) + 1
        Q(st,At) ← Q(st,At) +  $\frac{W}{C(st,At)} [G - Q(st,At)]$ 
        T(st) ← argmax Q(st,A)
        If At ≠ T(st) then exit for loop
        W ← W *  $\frac{1}{C(st,At)}$ 

```

⑧ TD(λ)

```

Initialize V(s) arbitrarily (ex. V(s)=0 for all s);
Repeat (for each episode):
    Initialize S
    Repeat (for each step of episode):
        At action given by π for s
        take action At, observe R,s'
        V(s) ← V(s) + δ(R + γV(s') - V(s))
        S ← s'
    until S is terminal

```

⑨ SARSA (On-policy TD control)

```

Initialize Q(s,a), Vs ∈ S, a ∈ A(s), arbitrary and Q(terminal-state)=0;
Repeat (for each episode):
    Initialize S
    choose A from S using policy derived from Q (ex. ε-greedy)
    Repeat (for each step of episode):
        take action A, observe R,s'
        choose A' from s' using policy derived from Q
        Q(s,A) ← Q(s,A) + δ(R + γQ(s',A') - Q(s,A))
        until S is terminal

```

⑩ Q-learning (Off-policy TD control)

```

Initialize Q(s,a), Vs ∈ S, a ∈ A(s), arbitrary and Q(terminal-state)=0;
Repeat (for each episode):
    Initialize S
    choose A from S using policy derived from Q (ex. ε-greedy)
    Repeat (for each step of episode):
        take action A, observe R,s'
        choose A' from s' using policy derived from Q
        Q(s,A) ← Q(s,A) + δ(R + γQ(s',A') - Q(s,A))
        until S is terminal

```

(11) Expected SARSA (On-policy / off-policy)

```

Q(s,At) ← Q(s,At) + δ(R + γ  $\sum_a \pi(a|s) Q(s',a) - Q(s,At)$ )
← Q(s,At) + δ(R + γ Q(s',At) - Q(s,At))

```

(12) Double Q-learning (off-policy)

```

Initialize Q1(s,a) and Q2(s,a), Vs ∈ S, a ∈ A(s), arbitrary
Initialize Q1(terminal-state) = Q2(terminal-state) = 0
Repeat (for each episode):
    Initialize S
    Repeat (for each step of episode):
        choose A from s using policy derived from Q1 and Q2 (ex. ε-greedy in Q1+Q2)
        Take action At, observe R,s'
        With probability:
            Q1(s,At) ← Q1(s,At) + δ(R + γ  $\arg\max_a Q_2(s',a) - Q_1(s,At)$ )
        else:
            Q2(s,At) ← Q2(s,At) + δ(R + γ  $\arg\max_a Q_1(s',a) - Q_2(s,At)$ )
        S ← s';
    until S is terminal

```

* <SARSA>

on-policy
choose an action (unecessary for best)
and update its value function with that knowledge
will converge eventually
slowly than learning

VS <Q-learning>

off-policy
choose an action
and update its value function with different action
no guarantee it will converge
faster & better than SARSA

Midterm

#1 The difference between MAB and full MDP is that in MABs, the full distribution over the results of actions are known T/F

The main difference is that MABs have no state. MDPs the distribution over results of an action are unknown.

#2 Exploration is when an agent is using its current estimate of the environment and value of the current state to take the action with highest reward even if it is not globally optimal T/F

The globally optimal action is always unknown w/o full info.

#3 To get the value function from a Q function, we need to get ~~mean~~ of the Q value of all ~~actions~~ at a given state. T/F
~~max state~~

#4 A coin is thrown, ball is drawn. If the coin is head, the ball is drawn from box 1 with 3R, 2B balls. If the coin is tail, Box 2 is used which has 2R, 3B balls.

a) $P(\text{red})?$ i) head: $\frac{1}{2} \times \frac{3}{5} = \frac{3}{10}$ ii) tail: $\frac{1}{2} \times \frac{2}{5} = \frac{1}{5} \therefore 0.4$

b) $P(\text{black})?$ i) head: $\frac{1}{2} \times \frac{2}{5} = \frac{2}{10} \therefore 0.2$

c) if the red ball is drawn, what is the probability that head is thrown? $\frac{0.3}{0.4} = 0.75$

d) $P(1 \text{ red})?$ $\frac{0.1}{0.4} = 0.25$

#5 Jane is 70% sure that her missing notebook is in one of 2 bags, being 40% certain it is in bag1 and 40% certain it is in bag2. If a search of bag1 doesn't reveal the notebook, what is the probability it is in bag2?

$$P(\text{bag2} | \text{not in bag1}) = \frac{P(\text{bag2, notebook}))}{P(\text{not in bag1})} = \frac{0.4}{1-0.4} = \frac{2}{3}$$

#6 Algorithm A is able to take input data D and provide a prediction of the most likely outcome of taking some action a in a particular state s.

Unfortunately, A is very unstable and doesn't converge easily. A modified algorithm B is proposed that improves stability by subtracting the current prediction (at time t) from the prediction from 3 steps earlier (at time t-3) and using the result to weight the value update. What kind of algorithm would you call B and explain?

Temporal difference method. Because algorithm B computes a difference between 2 time points to weight the value update.

#7 Let $G_t = R_t + \gamma R_{t+1} + \dots + \gamma^{T-t} R_T$. Suppose $r_{20:1}$ and the following sequence of rewards is received: $R_1=2, R_2=3, R_3=1, R_4=4$, $T=5$. Find G_1, G_2, G_3 .

$G_1 = \text{constant state}$

$$G_1 = R_1 + \gamma G_1 \Rightarrow 1 = 2 + \gamma G_1 \Rightarrow G_1 = -1$$

$$G_2 = R_2 + \gamma G_2 \Rightarrow -1 + \gamma G_2 = -1 + 0.1 \times 1 = -0.9$$

$$G_3 = R_3 + \gamma G_3 \Rightarrow 1 + 0.1 \times (-0.9) = 0.61$$

$$G_4 = R_4 + \gamma G_4 \Rightarrow 4 + 0.1 \times 0.61 = 4.061$$

$$\therefore G_5 = 9.061, G_6 = 1$$

#8 Algorithm Q-Learning

Algorithm parameters: stepsize $\alpha \in [0, 1]$, small $\epsilon > 0$ E number of episodes;

Initialise $Q(s,a)$ for all $s \in S^A, a \in A(s)$; arbitrarily except for $Q(\text{terminal}, \cdot) = 0$;

for each episode in E do:

 Report

$a' \leftarrow \arg\max Q(s,a)$;

 /* error : this below was actually on the exam.

$a' \leftarrow \arg\max Q(S_t, a_t)$;

 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 is terminal;

end

a) Is this algorithm an policy or off-policy? explain why

on-policy, since it is using the same action from $Q(s,a)$ to obtain both the behavioral action and the update action.

b) In what situations, if any would algorithm run into problems (learning)?

This algorithm could get stuck with suboptimal policies if it happens to find a feasible, but less than optimal path early on.

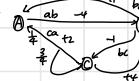
Since it does no random exploration and always maximizes, it will often stay with the first answer it finds that completes the task.

c) Assume now that this is meant to be an implementation of Q-learning, identify errors in the implementation of the above algorithm.

d) Should not be taken as the max of $Q(s,a)$, but rather the ϵ -greedy sample from $Q(s,a)$.

#9 Value function calculation.

Consider the following MDP with discount factor $\gamma = 0.9$. Upper case letter A,B,C : States, arc: state transitions, lower case letters ab,ba,ca,cb : actions, signed integers: rewards. Transitions: transition probabilities.



a) Consider the uniform random policy $\pi_1(s)$ that takes all actions from state s with equal probability. Starting with an initial value function of $V_1(A)=V_1(B)=V_1(C)=1$, apply one synchronous iteration of iterative policy evaluation (one backup for each state) to compute a new value function $V_2(s)$.

$$V_2(A) = \frac{1}{3} \pi_1(A) \sum_{s' \in S} P(s'|s,A) V_1(s') = \frac{1}{3}(4 + 2 + 2) = 2.67$$

$$V_2(B) = -4 + 0.9 V_1(C) = -4 + 0.9 \times 1 = -2.1$$

$$V_2(C) = \frac{1}{2} (4 + 0.9 V_1(A) + 0.9 V_1(B)) = 2 + 0.9 \times (2 + 0.9 \times 1) = 3.9$$

$$\therefore V_2(A) = 2.67, V_2(B) = 0.9, V_2(C) = 3.9$$

b) Apply one iteration of greedy policy improvement to compute a new, deterministic policy $\pi_2(s)$ assuming $V_2(A)=V_2(B)=2, V_2(C)=5$

$$\pi_2(s) = \arg\max_a \{ \pi_2(s, a) \mid \pi_2(s, a) = \frac{1}{|A(s)|} \sum_{a' \in A(s)} \pi_1(s, a') P(a'|s, a) V_2(s') \}$$

$$\pi_2(A, ba) = 1 + 0.9 V_2(C) = 1 + 0.9 \times 5 = 2.8$$

$$\pi_2(A, ab) = 1 + 0.9 V_2(C) = 1 + 0.9 \times 5 = 2.8$$

$$\pi_2(B, bc) = 2 + 0.9 \times (\frac{1}{2} V_2(A) + \frac{1}{2} V_2(C)) = 2 + 0.9 \times (\frac{1}{2} \times 2 + \frac{1}{2} \times 5) = 5.825$$

$$\pi_2(C, ca) = 2 + 0.9 V_2(A) = 2 + 0.9 \times 2 = 5.8$$

$$\pi_2(C, cb) = 2 + 0.9 V_2(A) = 2 + 0.9 \times 2 = 5.8$$

$$\therefore \pi_2(A) = ab, \pi_2(B) = bc, \pi_2(C) = ca$$

c) Apply one synchronous iteration (one backup for each state) to complete a new value function $V_2(s)$ assuming $V_1(A)=V_1(B)=V_1(C)=2$

$$V_2(A) = -4 + 0.9 V_1(C) = -4 + 0.9 \times 2 = -2.2$$

$$Q_2(B, ba) = 1 + 0.9 V_1(A) = 1 + 0.9 \times 2 = 2.8$$

$$Q_2(B, ab) = 1 + 0.9 V_1(C) = 1 + 0.9 \times 2 = 2.8$$

$$V_2(B) = \max(Q_2(B, ba), Q_2(B, ab)) = 2.8$$

$$Q_2(C, ca) = 2 + 0.9 V_1(A) = 2 + 0.9 \times 2 = 3.8$$

$$Q_2(C, cb) = 2 + 0.9 V_1(A) = 2 + 0.9 \times 2 = 3.8$$

$$V_2(C) = \max(Q_2(C, ca), Q_2(C, cb)) = 3.8$$

$$\therefore V_2(A) = -2.2, V_2(B) = 2.8, V_2(C) = 3.8$$

Assignment

1 MAB

Consider a k -armed Bandit problem with $k=4$, using ϵ -greedy action selection, sample average action-value

estimates, and initialization of $Q(a) = 0$, for all actions a . Suppose the initial sequences of actions & rewards: $A_1=2, A_2=1, A_3=4, A_4=4, R_1=0, R_2=2, R_3=4, R_4=3$

a) List out Q_t -estimates for each action, for each step from $t \in [5]$ before the action is taken. So row 1 is before any actions and row 6 is after the last action is taken

$$Q_t(a) = \frac{\sum_{i=1}^t R_i | A_i = a}{\sum_{i=1}^t 1 | A_i = a}$$

i) $t=0$ $Q_t(a)=0$

$$\text{ii) } t=2 \quad Q_2(a) = \frac{\sum_{i=1}^2 R_i | A_i = a}{\sum_{i=1}^2 1 | A_i = a}, \quad A_1=2 \quad \therefore Q_2(2)=0, \text{ o.w } Q_2(a)=0$$

$$\text{iii) } t=3 \quad Q_3(a) = \frac{\sum_{i=1}^3 R_i | A_i = a}{\sum_{i=1}^3 1 | A_i = a}, \quad A_1=2, A_2=1 \quad \therefore Q_3(1) = \frac{2}{3}, Q_3(2)=0, \text{ o.w } Q_3(a)=0$$

$$\text{iv) } t=4 \quad Q_4(a) = \frac{\sum_{i=1}^4 R_i | A_i = a}{\sum_{i=1}^4 1 | A_i = a}, \quad A_1=2, A_2=1, A_3=4 \quad \therefore Q_4(1) = \frac{2}{4} = 0.5, Q_4(2)=0, Q_4(4) = \frac{4}{4} = 1$$

$$\text{v) } t=5 \quad Q_5(a) = \frac{\sum_{i=1}^5 R_i | A_i = a}{\sum_{i=1}^5 1 | A_i = a}, \quad A_1=2, A_2=1, A_3=4, A_4=4 \quad \therefore Q_5(1)=2, Q_5(2)=0, Q_5(4)=\frac{2+4}{2}=3$$

$$\text{vi) } t=6 \quad Q_6(a) = \frac{\sum_{i=1}^6 R_i | A_i = a}{\sum_{i=1}^6 1 | A_i = a}, \quad A_1=2, A_2=1, A_3=4, A_4=4, A_5=2 \quad \therefore Q_6(2)=2, Q_6(4)=\frac{10}{5}=2, Q_6(1)=1.5, Q_6(3)=\frac{2+4}{2}=3$$

	t	$a=1$	$a=2$	$a=3$	$a=4$
Step 0	1	0	0	0	0
Step 1	2	0	0	0	0
Step 2	3	2	0	0	0
Step 3	4	2	0	0	2
Step 4	5	2	0	0	3
Step 5	6	2	1.5	0	3

b) Ecase: explore, greedy: exploit. Explore what we can tell about whether the agent definitely choose the ϵ case or whether it might have chosen the greedy case.

i) Step 0, all Q values = 0 for all the actions. \Rightarrow actions will be selected random $\therefore \epsilon$ action.

ii) Step 1, all Q values = 0 for all the actions. \therefore still ϵ action.

iii) $R_1=2$ for action 1 and other Q values are 0. So best action at this step is action 1.

But the agent selected action 4 $\therefore \epsilon$ action (definitely)

iv) maximum returns: both action 1 & 4. So one of it can be selected. Thus, may be ϵ action.

v) best action : 4, but it chose action 2 \therefore random (definitely)

c) Why can't we be any more sure than this?

The exploration ϵ case could have occurred at any step, even the cases where the greedy action is chosen, it could have been randomly selected, although this unlikely.

2 MDP

general form of the return for MDP from timestep t going forward until the terminal state is reached at $t=T$: $G_t = \sum_{k=t}^T r^{k-t} R_k = R_t + \gamma R_{t+1} + \dots + \gamma^{T-t} R_T$. But it also be defined recursively as for non-terminal states G_t :

Suppose $r=1$, and the following sequence of rewards: $R_1=6, R_2=3, R_3=0, R_4=-3, R_5=-6, T=5$. Once the terminal state is reached, there are no more rewards received and the episode ends.

a) For $G_0, G_1, G_2, G_3, G_4, G_5$, write out the one step recursive form and calculate.

$$G_0 = R_0 + 0.9 G_1 \rightarrow 6 + 0.9(-7.5) = -2.574$$

$$G_1 = R_1 + 0.9 G_2 \rightarrow 3 + 0.9(-7.5) = -7.56$$

$$G_2 = R_2 + 0.9 G_3 \rightarrow 0 + 0.9(-8) = -7.26$$

$$G_3 = R_3 + 0.9 G_4 \rightarrow -3 + 0.9(-6) = -8.7$$

$$G_4 = R_4 + 0.9 G_5 \rightarrow -6 + 0.9(-6) = -12$$

$$G_5 = R_5 + 0.9 G_6 \rightarrow 0 \text{ (terminal state)}$$

once terminal state reached, no reward $\therefore G_6=G_5$

$$\therefore G_0 = 2.574, G_1 = -7.56, G_2 = -7.26, G_3 = -8.7, G_4 = -12, G_5 = -12$$

b) calculate G_{10} by mnr recursive formula.

$$G_{10} = R_0 + \gamma R_1 + \gamma^2 R_2 + \gamma^3 R_3 + \gamma^4 R_4$$

$$= 6 + 0.9(3) + 0.9^2(0) + 0.9^3(-8) + 0.9^4(-6)$$

$$= 2.574$$

3 Policy evaluation

Figure shows a rectangular grid world representation of a simple finite MDP. The cells of the grid correspond to the states of the environment. At each cell, 4 actions are possible: north, south, east, west. Actions that would take the agent off the grid leave its location unchanged, but also result in a reward -1 . other actions result in a reward of 0 . With a discount factor $\gamma=0.9$, and equal probability

Specify Bellman equation for the setting, using the given equiprobable, random policy.

3.3	8.8	4.4	5.3	1.5
1.5	9.0	2.3	1.9	0.6
0.1	0.9	0.9	0.4	-0.4
-1	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2

a) Write down Bellman equation and describe what is known about each component in general terms.

$$V_{\pi}(s) = \sum_{a \in A(s)} \pi(a|s) P(s'|s, a) [r + \gamma V_{\pi}(s')].$$

For the given grid world, all actions have equal probability, so $\pi(a|s) = \frac{1}{4}$ for all actions, $r=0.9$.

b) Show the calculation for state J

$$\text{i) if it goes up: } \frac{1}{4}(0 + 0.9 \times 5.3) = 1.125$$

$$\text{ii) if it goes left: } \frac{1}{4}(0 + 0.9 \times 2.3) = 0.5135$$

$$\text{iii) if it goes right: } \frac{1}{4}(0 + 0.9 \times 0.5) = 0.1125$$

$$\text{iv) if it goes down: } \frac{1}{4}(0 + 0.9 \times 0.4) = 0.09$$

$$\therefore V_{\pi}(J) = 1.125 + 0.5135 + 0.1125 + 0.09 = 1.9125$$

$$\text{i) if it goes up: } \frac{1}{4}(0 + 0.9 \times 4.4) = -0.21$$

$$\text{ii) if it goes left: } \frac{1}{4}(0 + 0.9 \times 1.9) = -0.3125$$

$$\text{iii) if it goes right: } \frac{1}{4}(0 + 0.9 \times 1.3) = -0.21$$

$$\text{iv) if it goes down: } \frac{1}{4}(0 + 0.9 \times 1.2) = -0.21$$

$$\therefore V_{\pi}(K) = -0.21 - 0.3125 - 0.21 - 0.21 = -0.45 V_{\pi}(K)$$

$$V_{\pi}(K) = 1.79$$