



CMPUT365 Winter 2022

Final2: Sample Final

导师: Baihong Qi



Final Exam

1. (20)

For the average reward setting, the differential state value function is defined as:

$$v_{\pi}(s) \doteq E_{\pi} \left[\sum_{k=t+1}^{\infty} (R_k - r(\pi)) \mid S_t = s \right],$$

where the average reward is

$$r(\pi) = \lim_{h \rightarrow \infty} E_{\pi} \left[\frac{1}{h} \sum_{t=0}^{h-1} R_{t+1} \right].$$

Derive the following Bellman equation for v_{π} :

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{r, s'} p(s', r | s, a) [r - r(\pi) + v_{\pi}(s')].$$

Show steps.

LOTE: Law of total expectation

$$1. \quad E(X) = E(E(X | Y));$$

运算过程 [\[编辑\]](#)

$$\begin{aligned} E(E(X|Y)) &= \sum_y E(X|Y=y) \cdot P(Y=y) \\ &= \sum_y \left(\sum_x x \cdot P(X=x|Y=y) \right) \cdot P(Y=y) \\ &= \sum_y \sum_x x \cdot P(X=x|Y=y) \cdot P(Y=y) \\ &= \sum_y \sum_x x \cdot P(Y=y|X=x) \cdot P(X=x) \\ &= \sum_x \sum_y x \cdot P(Y=y|X=x) \cdot P(X=x) \\ &= \sum_x x \cdot P(X=x) \cdot \left(\sum_y P(Y=y|X=x) \right) \\ &= \sum_x x \cdot P(X=x) \\ &= E(X). \end{aligned}$$

MP: Markov Property

- Conditional Probability depends only upon the present states

LOTUS

Law of the unconscious statistician: $E[g(X)] = \sum_{x \in \mathcal{X}} g(x) P(X=x)$

Ans.

$$v_{\pi}(s) = E_{\pi} \left[\sum_{k=t+1}^{\infty} (R_k - r(\pi)) \mid S_t = s \right]$$

$$= E_{\pi} \left[R_{t+1} - r(\pi) + \sum_{k=t+2}^{\infty} (R_k - r(\pi)) \mid S_t = s \right]$$

$$= E_{\pi} \left[R_{t+1} - r(\pi) + E_{\pi} \left[\sum_{k=t+2}^{\infty} (R_k - r(\pi)) \mid S_{t+1}, S_t = s \right] \mid S_t = s \right]$$

$$= E_{\pi} \left[R_{t+1} - r(\pi) + E_{\pi} \left[\sum_{k=t+2}^{\infty} (R_k - r(\pi)) \mid S_{t+1} \right] \mid S_t = s \right] \quad \text{[LOTUS]}$$

$$= E_{\pi} \left[R_{t+1} - r(\pi) + v_{\pi}(S_{t+1}) \mid S_t = s \right] \quad \text{[M.P.]}$$

$$= \sum_{a, s', r} p_{\pi}(A_t = a, S_{t+1} = s', R_{t+1} = r \mid S_t = s) [r - r(\pi) + v_{\pi}(s')] \quad \text{[} v_{\pi} \text{ definition]}$$

$$= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + v_{\pi}(s')] \quad \text{[LOTUS]}$$

2. (20) Modify the Tabular TD(0) algorithm for estimating v_π , to estimate q_π .

Tabular TD(0) for estimating v_π

1. Input: the policy π to be evaluated
2. Algorithm parameter: step size $\alpha \in (0, 1]$
3. Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$
4. Loop for each episode:
 5. Initialize S
 6. Loop for each step of episode:
 7. $A \leftarrow$ action given by π for S
 8. Take action A , observe R, S'
 9. $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
 10. $S \leftarrow S'$
 11. until S is terminal

Sarsa

3. Sarsa算法(on-policy)

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

 until S is terminal

Ans.

* In step 3, replace " $V(s)$ " with " $Q(s, a)$ for all $a \in A$ "

and " $V(\text{terminal}) = 0$ " with " $Q(\text{terminal}, \cdot) = 0$ "

* Switch step 6 & 7

* Add before step 9 " $A' \leftarrow \text{action given by } \pi \text{ for } S'$ "

* Replace step 9 with

$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$$

* Add to step 10: " $A \leftarrow A'$ "

3. (20)

Show that the following off-policy TD(0) update for v_π is correct for transition $S, A \sim b, S', R$:

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

by showing the following:

$$\begin{aligned} & E_{A \sim b} [P(A|S)(R + \gamma V(S')) - V(S) | S=s] \\ &= E_{A \sim \pi} [R + \gamma V(S') - V(S) | S=s]. \end{aligned}$$

Here, $P(A|S) = \frac{\pi(A|S)}{b(A|S)}$, π and b are

policy distributions with assumption

$$\pi(A|S) > 0 \Rightarrow b(A|S) > 0.$$

Let's first work on the target

$$E_{A \sim b} [P(A|S)(R + \gamma V(S')) \mid S=s]$$

$$= \sum_{a, s', r} P_b(A=a, S'=s', R=r \mid S=s) P(a|s)(r + \gamma V(s'))$$

$$= \sum_{a, s', r} b(a|s) p(s', r \mid s, a) \frac{\pi(a|s)}{b(a|s)} (r + \gamma V(s'))$$

$$= \sum_{a, s', r} p(s', r \mid s, a) \pi(a|s) (r + \gamma V(s'))$$

$$= E_{A \sim \pi} [R + \gamma V(S') \mid S=s] \quad \textcircled{1}$$

On the other hand,

$$E_{A \sim b} [V(S) \mid S=s] = V(s).$$

$$\text{Therefore, } E_{A \sim b} [V(S) \mid S=s] = E_{A \sim \pi} [V(S) \mid S=s] \quad \textcircled{2}$$

Therefore, the identity follows from ①-②.

4. (20)

Give the specification of the off-policy
Expected Sarsa control method.

Sarsa

3. Sarsa算法(on-policy)

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

 until S is terminal

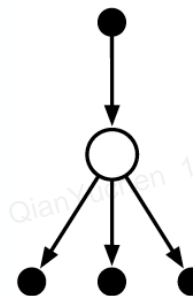
Sarsa

4. Expected Sarsa公式 (off-policy)

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

$$\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

- S_{t+1} 时look back
- current state 是 S_{t+1}



Expected Sarsa

5. Sarsa与Expected Sarsa的关系

- Sarsa是on-policy, Expected Sarsa多数情况下是off-policy

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

4. (20)

Give the specification of the off-policy
Expected Sarsa control method.

Ans. For the transition S, A, R, S' , where
the action A is drawn from policy
distribution b , update the action value
the following way:

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \sum_{a'} \pi(a'|S') Q(S', a') - Q(S, A)],$$

where the target policy π is off-policy and

uses greedification such as ϵ -greedy;

$\alpha > 0$, control

(20)

5. An agent is in a 3-state MDP, $S = \{1, 2, 3\}$, where each state has two actions $A = \{1, 2\}$. Assume the agent observed the following trajectory:

$$S_0 = 1, A_0 = 1, R_1 = 1, S_1 = 2, A_1 = 2, R_2 = -1,$$

$$S_2 = 3, A_2 = 1, R_3 = 2, S_3 = 1, A_3 = 1, R_4 = 2, S_4 = 1.$$

The agent uses Tabular Dyna-Q.

Which of the following are possible (or not possible) simulated transition $\{S, A, R, S'\}$ given the above observed trajectory with a deterministic model and random search control.

i. $\{S=1, A=1, R=2, S'=1\}$

ii. $\{S=3, A=1, R=2, S'=1\}$

iii. $\{S=1, A=1, R=1, S'=2\}$

iv. $\{S=2, A=2, R=-1, S'=3\}$

v. $\{S=3, A=2, R_3=2, S'=1\}$.

Just mention possible or not possible for each.

Planning & Learning

3. Dyna Q

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

(a) $S \leftarrow$ current (nonterminal) state

(b) $A \leftarrow \varepsilon$ -greedy(S, Q)

(c) Take action A ; observe resultant reward, R , and state, S'

(d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ ← **direct RL**

(e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) ← **model learning**

(f) Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ ← **planning**