Left page fragments:

...his proof, where there was

...he policy $\pi^*$, all inequalities

...have

...t factor $\alpha$) "solves"

$\sum P_{ij}(a)V_\alpha^*(j)]$

$a)V_\alpha^*(j)]$

reference state

$+ \alpha \sum P_{ij} h_\alpha(j)]$

better using this format

...ns, but

$\alpha \to 1$

---

**EE6106 LECTURE 23** Date: 17 Apr 2024

$\mu$: Classical RL $\to$ Q-learning, SARSA

   Discounted Cost Setting $(S, A, P, C, \alpha)$

Assumption : Underlying MDP is unichain. For average cost MDP, so such assumption was made when the horizon was assumed to be made.

In the RL setting, we require a connectivity assumption as mistakes are compulsory and sub-optimal actions have to be taken.

Flavours of RL :   1) On-policy learning
                   2) Off-policy learning

On-policy $\longrightarrow$ The learning algorithm is in control and decides which action to take.

Off-policy $\longrightarrow$ The entity which decides what action to take is different that the learning agent. We can still learn from the feedback signals received by someone else.

Flavours of RL algorithms : 1) Model-based
                            2) Model-free

Model-based $\longrightarrow$ Algorithm is learning a "model" and hence learning $(P, C)$

Model-free $\longrightarrow$ Not learning $P, C$ but directly trying to learn optimal policy. Learn value functions directly.

Note that we assume both $S$ and $A$ to be finite. Not doing stuff such as function approximation.

① Simple model-based on-policy approach

Idea : Play random actions until all $(s,a)$ pairs are seen

Then, Construct $\hat{P}_{ij}(a)$, $\hat{C}(i,a)$

$$\frac{\text{\# played } a \text{ and landed in } j \text{ from } i}{\text{\# played } a \text{ from } i}$$

For the policy, play $\epsilon$-greedy policy. i.e. w.p. $\epsilon_t$ play randomly & w.p. $1-\epsilon_t$ play optimal action for $(S,A,\hat{P},\hat{C},\alpha)$

$\epsilon_t$ is chosen to be a decreasing function of time.

Typically, $\epsilon_t \sim \frac{1}{t}$ or slower, such as $\epsilon_t \sim \frac{1}{\sqrt{t}}$. (Borel-Cantelli Lemma)

This ensures $\hat{P} \xrightarrow{a.s.} P$, $\hat{C} \xrightarrow{a.s.} C$ ← ensure infinite exploration w.p. 1

The memory footprint is $O(S^2 A)$

$\varepsilon_t$ is chosen to be a decreasing function of time.

Typically, $\varepsilon_t \sim \frac{1}{t}$ or slower, such as $\varepsilon_t \sim \frac{1}{\sqrt{t}}$. (Borel-Cantelli Lemma)

This ensures $\hat{P} \xrightarrow{a.s.} P$, $\hat{C} \xrightarrow{a.s.} C$ ⟵ ensure infinite exploration w.p. 1

The memory footprint is $O(S^2 A)$

↓ storing estimates of P.

[2] Q-Learning (model-free off-policy algorithm)

Learns the Q-function, with a memory requirement of $O(SA)$.

$$Q^*(i,a) = C(i,a) + \alpha \sum P_{ij}(a) V^*(j)$$
$$= C(i,a) + \alpha \sum P_{ij}(a)\left[\min_{a'} Q^*(j, a')\right]$$

This is the Bellman equation in terms of the Q function

$$\hat{Q}(x_t, A_t) \longleftarrow (1-\gamma_t)\hat{Q}(x_t, A_t) + \gamma_t\left[C_t + \alpha \min_{a'} \hat{Q}(x_{t+1}, a')\right]$$

$\gamma_t$ is the step-size parameter, gives memory of the earlier event

$\gamma_t$ is also chosen as a decreasing function of time

Note that $\min_{a'} \hat{Q}(x_{t+1}, a')$ is an estimate of $V^*(x_{t+1})$.

Need: $\sum \gamma_t = \infty$, $\sum \gamma_t^2 < \infty$

Q: What do we need for $\hat{Q}$ to converge to $Q$?

→ In the long run, every $(s,a)$-pair should be visited almost surely
Thus, we cannot have a stationary policy which drives Q-learning

"Tracks ODE" : $\dot{Q} = TQ - Q$, where $T$ is an operator

$$T u(i,a) = C(i,a) + \alpha \sum P_y(a)\left[\min_{a'} u(j,a')\right]$$

↓
a candidate
function

Note that the Bellman equation can be expressed as $TQ = Q$

The memory footprint now is $O(s,a)$.

③ SARSA (State, Action, Reward, State, Action)
Model-free, on-policy algorithm

$$A_t = \begin{cases} \text{Play random action w.p. } \varepsilon_t \\ \\ \arg\min_a \hat{Q}(x_t, a) \text{ w.p. } 1 - \varepsilon_t \end{cases}$$

$$\hat{Q}(x_t, A_t) \leftarrow (1-\gamma_t)\hat{Q}(x_t, A_t) + \gamma_t\left[C_t + \alpha\hat{Q}(x_{t+1}, A_{t+1})\right]$$

There is no minimization term in SARSA. Note that the next step
we WILL perform minimization w.p. $1 - \varepsilon_t$
action selection

$$A_t = \begin{cases} \text{Play random action w.p.} \\ \text{avg min } \hat{Q}(x_t, a) \quad \text{w.p. } 1 - \varepsilon_t \\ \quad a \end{cases}$$

$$\hat{Q}(x_t, A_t) \leftarrow (1 - \gamma_t) \hat{Q}(x_t, A_t) + \gamma_t \left[ C_t + \alpha \hat{Q}(x_{t+1}, A_{t+1}) \right]$$

There is no minimization term in SARSA. Note that the next step we WILL perform minimization w.p. $1 - \varepsilon_t$

Note that for $Q$-learning, we can use any action selection algorithm as long as it performs infinite exploration.

Typically, choose $\varepsilon_t \sim \frac{1}{t}$ & $\gamma_t \sim \frac{1}{t}$

$$\Rightarrow \hat{Q} \xrightarrow{\text{a.s.}} Q$$

Note that these algorithms are notoriously slow. This is because only entry gets filled in one time-step

Function approximation : expresses value function as a linear combination of some basic functions.

---

4 Actor - Critic algorithm

GPI $\longrightarrow$ Generalized Policy Iteration

Original policy iteration searches for greedy improvements on the policy

Actor : Policy improvement   (slower timescale)
Critic : Policy evaluation   (faster timescale)

You have two timescales, one with a small & one with a large step size

In practice, use   step size $\gamma_t$ for actor & $\beta_t$ for critic.

$$\beta_t = o(\gamma_t) \quad \gamma_t = o(\beta_t)$$

Little of Oh

## SARSA

If we do, off-policy SARSA, under some stationary policy $\pi$, we will learn the value function corresponding to the policy $\pi \longrightarrow Q^\pi$ and $V^\pi$

All of these algorithms can be generalized for a time-average MDP setting.

Lecture 23

Page No.