

Lecture 19

Monday, October 26, 2020

CS 59000-RL

MDP

- Stochastic approximation
- Q-learning
- Infinite horizon MDP (undiscounted)

for off policy vs on policy } check out
Q-learning SARSA Rich, Andrew's
Sutton Berbo
An Intro to RL

Let's see when Q-learning may converge:
we use tools from stochastic approximation Theory
Robbins-Munro 1951, Dvoretzky 1956

Theorem (Jaakkola et al. 1993)

A random iterative process

$$\rightarrow \Delta_{t+1}(x) = (1 - \alpha_t(x)) \Delta_t(x) + \beta_t(x) F_t(x)$$

converges to zero a.s. if for all $x \in \Omega$

the following condition holds:

1) The Ω is finite.

$$2) \sum_t \alpha_t(x) = \infty, \sum_t \alpha_t^2(x) < \infty$$

$$\sum_t \beta_t(x) = \infty, \sum_t \beta_t^2(x) < \infty$$

$$E[\beta_t(x) | \mathcal{F}_t] \leq E[\alpha_t(x) | \mathcal{F}_t] \text{ a.s.}$$

$$3) E[\underbrace{F_t(x)}_{\text{}} | \mathcal{F}_t] \leq \gamma \cdot E[\|\Delta_t\| | \mathcal{F}_t] \quad (0 < \gamma < 1)$$

$$4) \text{Var}[F_t(x) | \mathcal{F}_t] \leq C \left(1 + E[\|\Delta_t\| | \mathcal{F}_t]\right)^2$$

for some constant C .

where $\{\Delta_t, \Delta_{t-1}, \dots, F_{t-1}, F_{t-2}, \dots, \alpha_{t-1}, \alpha_{t-2}, \dots, \beta_{t-1}, \beta_{t-2}, \dots\}$
is \mathcal{F}_t measurable.

Theorem: The Q learning algorithm given by

$$Q_{t+1}(x, a) = (1 - \alpha_t(x, a)) Q_t(x, a) + \alpha_t(x, a) \left[r(x, a) + \gamma \max_{a'} Q_t(x', a) \right]$$

$$r(x, a) \sim R(x, a)$$

$$x' \sim P(\cdot | x, a)$$

Converges to the optimal Q^* if:

- 1) The state and action spaces are finite.
- 2) $\sum_b \alpha_b(x, a) = \infty$ and $\sum_b \alpha_b^2(x, a) < \infty$ a.s.
- 3) $\text{Var}[r(x, a) | x, a]$ is bounded.

Proof: Let set $\Delta_t(x, a) = Q_t(x, a) - Q^*(x, a)$

$$\rightarrow F_t(x, a) = r_t(x, a) + \gamma \max_{a'} Q_t(x', a') - Q^*(x, a)$$

$$\beta_t = \alpha_t$$

$$\downarrow$$

$$E[r(x, a) + \max_{a'} Q^*(x', a')]$$

— $X \times A = \Omega$ is finite. ☺

↳ The first condition

— Since $\alpha_t \geq \beta_t \Rightarrow$ The second condition is satisfied. ☺

— For the third one.

$$\begin{aligned} \max_a \left| E[F_t(x, a) | \mathcal{F}_t] \right| &= \max_a \left| E \left[\cancel{r_t(x, a)} + \gamma \max_{a'} Q_t(x', a') - \cancel{r_t(x, a) - \gamma \max_{a'} Q^*(x', a')} \right] \right| \\ &= \gamma \max_a \left| \sum_{x'} P(x' | x, a) \left(\max_{a'} Q_t(x', a') - \max_{a'} Q^*(x', a') \right) \right| \end{aligned}$$

Monday, October 26, 2020

$$\leq \gamma \max_a \sum_{n'} p(n'|n, a) \max_{a'} |Q_b(n', a') - Q^*(n', a')|$$

$$\|Q_t - Q^*\|_\infty$$

$$\leq \gamma \|Q_t - Q^*\|_\infty = \gamma \|\Delta_t\|_\infty$$

we showed

$$\|E[F_t(\cdot) | \mathcal{F}_t]\|_\infty = \max_{n, a} \|E[F_t(n, a) | \mathcal{F}_t]\|_\infty$$

$$\leq \gamma \|\Delta_t\|_\infty$$

Now the force one:

$$\text{var}[F_t(n, a) | \mathcal{F}_t] =$$

$$E\left[\left(\underbrace{r_t(n, a)}_{\text{red}} + \gamma \max_{a'} Q_b(n', a') - Q^*(n, a)\right)^2 \middle| \mathcal{F}_t\right]$$

$$= E\left[\left(\underbrace{r_t(n, a) - \bar{r}(n, a)}_{\text{red}} + \gamma \sum_{n'} p(n'|n, a) \max_{a'} Q_b(n', a') - Q^*(n, a)\right)^2 \middle| \mathcal{F}_t\right]$$

$$= E\left[\left(r_t(n, a) - \bar{r}(n, a)\right)^2 \middle| \mathcal{F}_t\right]$$

$$+ \gamma^2 E\left[\left(\max_{a'} Q_b(n', a') - \sum_{n'} p(n'|n, a) \max_{a'} Q_b(n', a')\right)^2 \middle| \mathcal{F}_t\right]$$

Monday, October 26, 2020

$$\leq \text{Var}[r(x, a)] + \gamma^2 E\left[\left(\max_{a'} Q_t(x', a')\right)^2\right]$$

$$\leq \text{Var}[r(x, a)] + \gamma^2 E\left[\max_{a'} (Q^*(x', a'))^2\right] + \gamma^2 \|\Delta_t\|_\infty^2$$

$$\leq \max_{x, a} \text{Var}[r(x, a)] + \gamma^2 E\left[\max_{a'} (Q^*(x', a'))^2\right] + \gamma^2 \|\Delta_t\|_\infty^2$$
$$\leq C(1 + \|\Delta_t\|_\infty^2)$$

which satisfies the fourth condition.

In finite horizon undiscounted MDP

$$M: (X, A, P, P_1, R)$$

Expected return under a policy π

$$E^\pi\left[\sum_{t \geq 0} r_t\right] \rightarrow \text{expected cumulative reward.}$$

$\sum_{t \geq 0} r_t$ if the sum exists, might not be integrable
- But some time we like.

For strongly connected MDPs, define long term average expected return of a memory-less π starting from:

$$\rho_\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E^\pi[r(x_t, A_t) | X_1 = x]$$

(Assumption, the limit exists)

Define $\rho^\pi = \max_x \rho_x^\pi$, since M is strongly connected and π is memoryless, for optimal policy

$$\rho^* = \max_{\pi \rightarrow \text{memoryless}} \bar{\rho}^\pi$$

$$\rho^* = \rho_x^{\pi^*} \text{ for } x \in \mathcal{X}$$

Back to matrix notation.

$$\rho^\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T P_\pi^{t-1} r_\pi$$

Define $\bar{P}_\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T P_\pi^{t-1}$ as the stationary transition matrix, $\Rightarrow \rho^\pi = \bar{P}_\pi r_\pi$

Starting from a state $X_1 = x$ does not change the average expected return, but has some gain.

Value function:

first define $V_\pi^{(T)} = \sum_{t=1}^T P_\pi^{t-1} (r_\pi - \rho^\pi)$

The value function is the Cesaro sum of $P_\pi^t (r_\pi - \rho^\pi)$

$$V_\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T V_\pi^t \leftarrow$$

Monday, October 26, 2020

$V_{\pi}(x) - V_{\pi}(x')$ represents long term advantage
of observing from x v.s. x' .

$$\text{Lemma: } V_{\pi} = \left((I - P_{\pi} - \bar{P}_{\pi})^{-1} - \bar{P}_{\pi} \right) r_{\pi}$$

proof: HW.

Monday, October 26, 2020

Monday, October 26, 2020

Monday, October 26, 2020

Monday, October 26, 2020

Monday, October 26, 2020

Monday, October 26, 2020