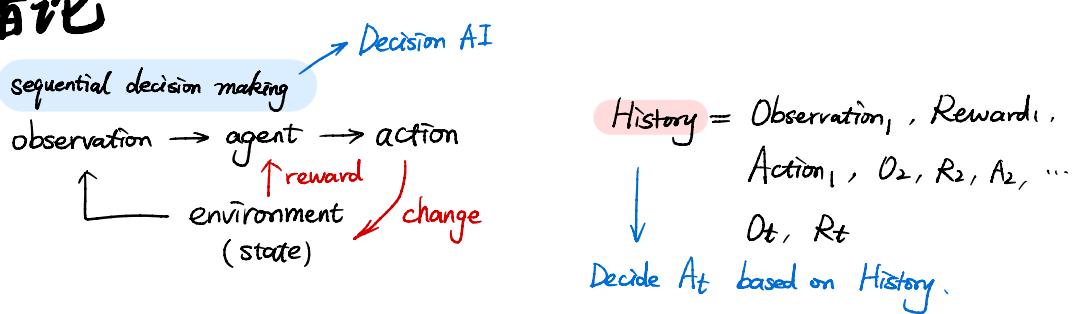


强化学习

Lecturer: 李少卿

Note Taker: EP-范元

绪论



State: 描述环境的所有信息，用于确定将发生的事 (O_t, A_t, R_t)

Assumption. $S_t = f(H_t)$. 环境可以从 H_t 中判断出 S_t .

policy: $\pi: S \rightarrow \Delta(A)$

$$\begin{cases} \text{deterministic} & a = \pi(s) \\ \text{stochastic} & \pi(a|s) = \Pr[A_t=a | S_t=s] \end{cases}$$

reward: $R(s, a, s')$ ~ 感知到什么是“好”的 (for long term)

utility of a sequence.

$$[r_0, r_1, \dots] \quad \begin{matrix} \text{more vs less} \\ \text{now vs later} \end{matrix} \rightarrow \text{discounting } \gamma. \quad r_0 + \gamma r_1 + \gamma^2 r_2 + \dots = \text{Utility}$$

[e.g. interest]

Theorem. $[a_1, a_2, \dots] \succ [b_1, b_2, \dots] \Leftrightarrow [r, a_1, a_2, \dots] \succ [r, b_1, b_2, \dots]$

value function predict reward in the future → decide whether the policy is good or not.

$$\begin{aligned} Q_{\pi}(s, a) &= \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t=s, a_t=a] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma Q_{\pi}(s', a') | S_t=s, a_t=a] \end{aligned}$$

model predict the next state

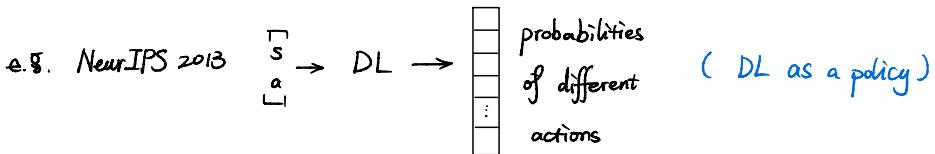
$$P_{ss'}^a = \Pr[S_{t+1}=s' | S_t=s, A_t=a] \quad \text{转移概率}$$

$$R_s^a = \mathbb{E}[R_{t+1} | S_t=s, A_t=a] \quad \text{预测 reward (一些问题中是确定且已知的)}$$

$$\left\{ \begin{array}{l} \text{Based on value: } a^* = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a) \rightarrow \text{policy: hidden} \\ \text{Based on policy: explicit form, without value} \\ \text{Actor-Critic 学习 / 模仿} \end{array} \right.$$

RL

$(s, a) \rightarrow Q(s, a)$ large-scale!
近似两者之间的函数 $\sim DL$. } Deep RL. (DL as value function)



Actually what we want of Q is

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a) = \mathbb{E}[r_{t+1} + \gamma Q^*] = \mathbb{E}[r_{t+1} + \gamma \max_{\pi} Q_{\pi}]$$

Q-Learning

$$Q_{\theta}(s_t, a_t) \leftarrow (1-\alpha) Q_{\theta}(s_t, a_t) + \alpha (r_t + \gamma \max_{a'} Q_{\theta}(s_{t+1}, a'))$$

$\xleftarrow{\text{参数}} Q_{\theta}(s_t, a_t) + \alpha (r_t + \gamma \max_{a'} Q_{\theta}(s_{t+1}, a') - Q_{\theta}(s_t, a))$

$$L(\theta) = (Q_{\theta} - \underbrace{(r + \gamma Q)}_{\text{target. 也包括 } Q})^2$$

$\left\{ \begin{array}{l} Q_{\theta}: \text{evaluation network} \\ Q: \text{target network, 是 } Q_{\theta} \text{ 几轮前的结果. 更新慢一些} \end{array} \right.$

前沿

<p>simulator. 基于模拟模型的强化学习</p> <p>goal-oriented <robotics> ~ 长期任务挑战 → 分割成短期简单任务 (中间态)</p> <p>模仿学习 (Imitation Learning)</p> <p>多智能体强化学习 ~ 非稳态环境 How to collaborate?</p>

exploitation v.s. exploration

π : 当前策略. $\xrightarrow{\text{保持已知的最优}} \text{(exploitation)}$
 $\xrightarrow{\text{探索更多决策, 但未必最优.}} \text{(exploration)}$

多臂老虎机 $A = \{a^1, a^2, \dots, a^K\}$. $R(r, a^i) = \Pr[r | a^i]$

AIM: $\max \sum_{t=1}^T r_t$. $r_t \sim R(\cdot, a_t)$.

$$Q_n(a^i) = \frac{r_1 + \dots + r_{n-1}}{n-1} \quad Q_{n+1}(a^i) = Q_n(a^i) + \frac{1}{n} (r_n - Q_n)$$

regret

$$R(a^i) = Q^* - Q(a^i) = \max_{a^j \in A} Q(a^j) - Q(a^i)$$

$$Q(a^i) = \mathbb{E}_{r \sim \Pr[r | a^i]} (r)$$

total regret $\Omega_R = \mathbb{E}_{a \sim \pi} \left[\sum_{t=1}^T R(a_t^i) \right]$

$$\min \sigma_R \Leftrightarrow \max \mathbb{E}_{a \sim \pi} \left[\sum_{t=1}^T Q(a_t^i) \right] \quad \text{-直 exploit / -直 explore} \sim \text{线性递增}$$

$$\sigma_R = \mathbb{E} \left[\sum_{t=1}^T Q(a^*) - \sum_{t=1}^T Q(a_t) \right] \quad \frac{\sigma_R}{T} \rightarrow 0 \Leftrightarrow Q(a_t) \rightarrow Q(a^*)$$

$$\sigma_R(T) = o(T), \quad \text{sublinear}$$

$$\lim_{T \rightarrow \infty} \sigma_R \geq \log T \sum_{a | \Delta_a > 0} \frac{\Delta_a}{D_{KL}(R(r|a) \| R^*(r|a))} \sim \log T \sum_{a | \Delta_a > 0} \frac{1}{\Delta_a}. \quad \begin{array}{l} \text{explore} \\ \log T \text{ not optimal} \\ T - \log T \text{ optimal} \end{array}$$

greedy: 看看目前数据里哪个 arm 最好, 用它. $\Rightarrow \varepsilon\text{-greedy}$: 以 ε 概率选择 $U(1, N)$. \downarrow 衰减 ε 能指数下降

乐观初始化 reward $\in (0, 1)$. \rightarrow 初始化: $Q(a^i) >> 1$, 如用 5 或 10.

增量更新 $Q(a^i)$, 从 5 降到 $(0, 1)$ 范围. \rightarrow 未 pull 的 arm 保持很高的 $Q(\cdot)$.

\rightarrow 各个 arm 都連續 pull 多次才能回到较正常的估值

$$\min \left\{ 1, \frac{c \log t}{t^{2-\delta}} \right\}$$

\uparrow
need knowledge

\Rightarrow 每个 arm 都被 pull 了较多次数 (exploration).

动作的价值分布



如何平衡 exploitation vs exploration

不确定性大的 $Q(a^i)$ 有探索价值

UCB Hoeffding 不等式. $\Pr[\mathbb{E}[x] > \bar{x}_t + u] \leq e^{-2tu^2}$ for $x \in [0, 1]$

估计上置信界 $\hat{U}(a^i) \Rightarrow Q_t(a) \leq \hat{Q}_t(a) + \hat{U}_t(a)$ 以高概率 $(1-p)$ 成立

$$a^* = \underset{a \in A}{\operatorname{argmax}} \hat{Q}_t(a) + \hat{U}_t(a). \quad e^{-2N_t(a)U_t(a)^2} = p \Rightarrow \hat{U}_t(a) = \sqrt{-\frac{\log p}{2N_t(a)}} \quad \text{Hoeffding}$$

\downarrow exploration

Thompson Sampling 根据每个动作成为最优动作的概率来选择. 根据概率分布采样, 选择 $\operatorname{argmax} \tilde{Q}(\cdot)$ \leftarrow 采样.

each arm \rightarrow Init: Gaussian $(0, 1)$

$$\text{posterior} \Rightarrow \text{Gaussian} \left(\hat{\alpha}_a(t), \frac{1}{1 + T_a(t)} \right) \quad \leftarrow \text{采样越多越确信}$$

Markov Decision Process

Markov 性

the future is independent of the past when given the present.

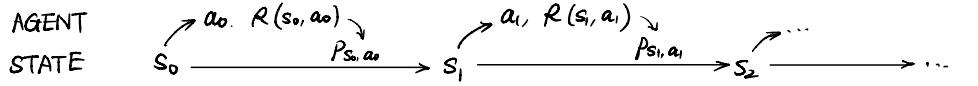
$$\Pr[S_{t+1} | S_t] = \Pr[S_{t+1} | S_1, \dots, S_t]$$

$$\Pr[S_{t+1} | S_t, A_t]$$

状态从历史中获得了充分信息
不管历史 (S_t 吸收了 H_t 中所有信息)

$$\pi: H_t \rightarrow A \quad \xrightarrow{\text{Markov Property}} \quad \pi: S \rightarrow A$$

MDP: $(S, A, \{P_{s,a}\}, \gamma, R)$. for RL: P_{sa} and R : unknown



TOTAL REWARD. $R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots$ 多数情况下 R 只与 s 有关.

Given an environment (MDP in this case), different strategies samples a different distribution of (s, a) .

occupancy measure $\rho^\pi(s, a) = \sum_{t=0}^T \gamma^t \Pr[S_t = s, a_t = a | S_0, \pi]$

$$\begin{aligned} \rho^\pi(s) &= \sum_{t=0}^T \gamma^t \Pr[S_t = s | S_0, \pi] = \sum_{t=0}^T \gamma^t \Pr[S_t = s | S_0, \pi] \sum_{a'} \pi(a_t = a' | S_t = s) \\ &= \sum_a \sum_{t=0}^T \gamma^t \Pr[S_t = s, a_t = a | S_0, \pi] = \sum_a \rho^\pi(s, a) \end{aligned}$$

Thms. $\rho^{\pi_1} = \rho^{\pi_2}$ iff. $\pi_1 = \pi_2$. $\pi_\rho(a|s) = \frac{\rho(s, a)}{\sum_{a'} \rho(s, a')}$

TOTAL REWARD. $V(\pi) = \mathbb{E}_\pi[R(s, a)]$

也可以直接解 $V^\pi = r + \gamma \cdot P_V^\pi$

已知 $P_{s,a}$ 和 R 求解策略

$$V^\pi(s) = \mathbb{E}_\pi[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | S_0 = s]$$

$$= R(s) + \gamma \sum_{s' \in S} P_{s, \pi(s)}(s') V^\pi(s') = R(s) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s')$$

value iteration

要求 $v_0 = f(v_0)$ 是 γ -Lipschitz (contracting). i.e. $\|f(v_0) - f(v'_0)\| \leq \gamma \|v_0 - v'_0\|$.

$$\begin{aligned} v_0 &\rightarrow v^{(1)} = f(v_0) \rightarrow \dots \rightarrow v^{(n)} = f(v^{(n-1)}) \\ v'_0 &\rightarrow v^{(1)'} = f(v'_0) \rightarrow \dots \rightarrow v^{(n)'} = f(v^{(n-1)'}) \end{aligned} \quad \left\{ \begin{array}{l} \|v^{(n)} - v^{(n)'}\| \leq \gamma^n \|v_0 - v'_0\| \\ (\text{收敛到同一点}) \end{array} \right. \rightarrow 0.$$

Bellman Equation $V^*(s) = \max_\pi V^\pi(s) \Rightarrow V^*(\pi) = R(s) + \gamma \max_{a \in A} \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\pi(s)$ (still contracting.)

policy iteration

$$\pi'(s) = \arg\max_a R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\pi(s) \rightsquigarrow \text{Is } V^{\pi'} \geq V^\pi? \quad \begin{cases} \pi = \pi^*, & V^{\pi'} = V^\pi \\ \pi \neq \pi^*, & V^{\pi'} > V^\pi \end{cases}$$

(converge faster)

* 从经验中学习 MDP 模型 → 动态规划 (v/p Mer)

MODEL-FREE RL

(model-based)

用数据逼近概率

通常不会给出状态转移和 reward \rightarrow ① 估计模型 MDP ② Sampling (model-free)

Monte-Carlo

$$V^\pi(s) \approx \frac{1}{N} \sum_{i=1}^N G_t^{(i)}$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

增量 MC 更新

$$N(s_t) += 1$$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t))$$

← 从完整的片段中学习

(只适用于有限长度 MDP)

对非稳定问题 \sim 追踪现阶段平均值

以 α 来减

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$$

$$\mathbb{E}_{\text{exp}}[f(x)] = \int_x p(x) f(x) dx = \int_x g(x) \frac{p(x)}{g(x)} f(x) dx = \mathbb{E}_{x \sim g} \left[\frac{p(x)}{g(x)} f(x) \right]$$

评估：使用 μ 策略产生的奖励评估 π 策略

$$\text{重要性比率: } \frac{\pi(s)}{\mu(s)}$$

$$\{s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T\} \sim \mu \Rightarrow G_t^{\pi/\mu} = \frac{\pi(a_1|s_1)}{\mu(a_1|s_1)} \frac{\pi(a_2|s_2)}{\mu(a_2|s_2)} \dots \frac{\pi(a_T|s_T)}{\mu(a_T|s_T)} G_t$$

使用重要性采样的离线 MC

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^{\pi/\mu} - V(s_t))$$

时序差分学习

(TD)

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = R_{t+1} + \gamma V(s_{t+1})$$

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

observation

prediction for future

← 直接从经验片段中学习

$$\text{sample} = R(s, \pi(s), s') + \gamma V^\pi(s') \Rightarrow V^\pi(s) \leftarrow (1-\alpha) V^\pi(s) + \alpha \cdot \text{sample} \sim \text{sample} \Rightarrow V^\pi \text{ 的无偏估计.}$$

① Gradient Descent View

$$\text{i.e. } V^\pi(s) \leftarrow V^\pi(s) + \alpha \cdot \underline{[\text{sample} - V^\pi(s)]}$$

$$\text{error}(V) = \frac{1}{2} (\text{sample} - V^\pi(s))^2 \rightarrow \nabla \text{error}.$$

② Exponential Moving Average View

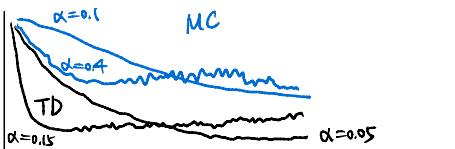
$$\sim V_n = \alpha x_n + \alpha(1-\alpha) x_{n-1} + \dots + \alpha(1-\alpha)^{n-2} x_2 + (1-\alpha)^{n-1} x_1$$

"forgets the past".

$$\alpha = \frac{1}{N}$$

MC $\alpha \uparrow \sim$ new sample would greatly influence V . e.g. $\alpha = 0.04 \rightarrow$ 到 25 之前快速收敛 (光滑)
TD $\alpha = 0.15$

25 之后剧烈振荡
(α 过大, new sample 的方差影响大)



多步时序差分学习

n 步累计奖励 $G_t^{(n)} = R_{t+1} + \gamma \cdot R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t))$$

γ^n : 使 bias 减少 \rightarrow bias 的主要来源.

How to improve policies?



$P(s, a, s')$: unknown! \Rightarrow Impossible to get $\mathcal{Q}(s, a)$

$\pi(s) \rightarrow$ we need to observe at state s . \leftarrow Fix s ? Hard.

MC and TD can't.

(开发者模式)(x)

SARSA

\sim learn $\mathcal{Q}^\pi(\cdot)$. $\mathcal{Q}^\pi(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \cdot \mathcal{Q}^\pi(s', \pi(s'))]$

$$\mathcal{Q}(s, a) \leftarrow \mathcal{Q}(s, a) + \alpha (r + \gamma \mathcal{Q}(s', a') - \mathcal{Q}(s, a)).$$

improvement of π : ϵ -greedy $\begin{cases} \underset{a}{\operatorname{argmax}} \mathcal{Q}(s, a) & \text{w.p. } 1-\epsilon \\ U[a] & \text{w.p. } \epsilon. \end{cases}$

on-policy. \sim 使用当前策略进行动作采样. (基于当前策略选择下一个动作) [实时改善].

Q-Learning

AIM: Learn \mathcal{Q}^* . $\mathcal{Q}^*(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \cdot \max_{a'} \mathcal{Q}^*(s', a')]$

{ 目标策略 π \rightarrow 用以评估 V^π . \mathcal{Q}^π

{ 行为策略 μ \rightarrow 用以收集数据. $\{s_1, a_1, r_2, s_2, a_2, \dots, s_T\} \sim \mu$.
(在环境中尝试)

环境 $a_t \sim \mu(\cdot | s_t)$.

offline $a'_{t+1} \sim \pi(\cdot | s_t)$. $\mathcal{Q}^*(s_t, a_t) \leftarrow \mathcal{Q}^\pi(s_t, a_t) + \alpha (r_{t+1} + \gamma \mathcal{Q}^\pi(s_{t+1}, a'_{t+1}) - \mathcal{Q}(s_t, a_t))$

用 μ 在环境中撞出来的 $s_t, a_t \rightarrow s_{t+1}$ 来更新 \mathcal{Q}^* .

π : $\pi(s_{t+1}) = \underset{a'}{\operatorname{argmax}} \mathcal{Q}^*(s_{t+1}, a')$ } $r_{t+1} + \gamma \mathcal{Q}(s_{t+1}, a'_{t+1}) = r_{t+1} + \gamma \mathcal{Q}(s_{t+1}, \underset{a'_{t+1}}{\operatorname{argmax}} \mathcal{Q}(s_{t+1}, a'_{t+1}))$

μ : ϵ -greedy policy of $\mathcal{Q}^*(s, a)$.

$$= r_{t+1} + \gamma \max_{a'_{t+1}} \mathcal{Q}(s_{t+1}, a'_{t+1}).$$

$$\mathcal{Q}(s_t, a_t) \leftarrow \mathcal{Q}(s_t, a_t) + \alpha (r_{t+1} + \gamma \max_{a'_{t+1}} \mathcal{Q}(s_{t+1}, a'_{t+1}) - \mathcal{Q}(s_t, a_t))$$

用以学到 \mathcal{Q}^* . (而非 \mathcal{Q}^π).
 $\mathcal{Q} \rightarrow \mathcal{Q}^*$ eventually.

完全反向传播(DP) 全展开

采样反向传播(TD)

$V^\pi(s)$

value iteration

时序差分

$Q^\pi(s, a)$

Q -policy iteration

SARSA

$Q^*(s, a)$

Q -value iteration

Q -Learning

多步时序差分学习(reprise)

$$G_t^{(n)} = R_{t+1} + \gamma \cdot R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}).$$

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t)) \quad \alpha \uparrow \sim \text{forget more about past. 新采样影响更大.}$$

$$\sum \lambda^n = \frac{1}{1-\lambda}$$

\downarrow \$\left\{ \begin{array}{l} n \uparrow \cdots \text{Monte Carlo. (能承受的 } \alpha \text{ 更小)} \\ n \downarrow \cdots \text{TD (能承受的 } \alpha \text{ 更大)} \end{array} \right.
用以控制一下方差

$$\Rightarrow G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \sim \text{结合平均不同步数下的奖励.} \quad \left\{ \begin{array}{l} \text{...} \\ \text{...} \\ \text{...} \\ \text{...} \\ \text{...} \end{array} \right. \quad \left\{ \begin{array}{l} 1-\lambda \\ (1-\lambda)\lambda \\ (1-\lambda)\lambda^2 \\ (1-\lambda)\lambda^3 \\ \dots \end{array} \right.$$

$$\lambda = 1 \rightarrow \text{Monte Carlo.} \quad \lambda = 0 \rightarrow \text{TD}$$