

RL 基础

集中不等式：

① Markov Inequality

$X \in \mathbb{R}^+$, 有 $\Pr(X \geq s) \leq \frac{E(X)}{s}$

$$\text{Prof: } E(X) = \int_0^s x p(x) dx + \int_s^\infty x p(x) dx$$

$$\geq \int_s^\infty x p(x) dx$$

$$\geq s \int_s^\infty p(x) dx = s \Pr(X \geq s)$$

$$\Rightarrow \Pr(X \geq s) \leq \frac{E(X)}{s}$$

② Chebyshev's Inequality

$$\Pr(|X - E(X)| \geq \epsilon) = \Pr((X - E(X))^2 \geq \epsilon^2)$$

$$\leq \frac{E[(X - E(X))^2]}{\epsilon^2} = \frac{\text{VAR}(X)}{\epsilon^2}$$

③ Chernoff Lemma \rightarrow Chernoff Bound

$$\text{Prof: } \Pr(X - E(X) \geq \epsilon) \leq \min_{\lambda \geq 0} \frac{E[\exp(\lambda(X - E(X)))]}{\exp(\lambda\epsilon)} \quad \text{and} \quad \Pr(X \geq \epsilon) \leq \min_{\lambda \geq 0} \frac{E[\exp(\lambda X)]}{\exp(\lambda\epsilon)}$$

$$\text{Prof: } \Pr(X - E(X) \geq \epsilon) \stackrel{\lambda \geq 0}{=} \Pr(\exp(\lambda(X - E(X))) \geq \exp(\lambda\epsilon)) \leq \frac{E[\exp(\lambda(X - E(X)))]}{\exp(\lambda\epsilon)} \rightarrow \text{MGF} \rightarrow \frac{E[\exp(\lambda X)]}{\exp(\lambda\epsilon)} = \frac{E[\exp(\lambda X)]}{\exp(\lambda\epsilon + \lambda\epsilon)}$$

$$\text{Chernoff Inequality} \quad X_i \sim P_x(\cdot) \quad \Pr\left(\sum_{i=1}^n X_i - nE[X] \geq \epsilon\right) \leq \min_{\lambda \geq 0} \left[\prod_{i=1}^n (E[\exp^{\lambda(X_i - E[X))}]) \right] \exp^{-\lambda\epsilon}$$

$$\text{Prof: } \text{MGF}_{X_1+X_2+\dots+X_n}(\lambda) = \prod_{i=1}^n \text{MGF}_{X_i}(\lambda)$$

Chernoff Bound: Get the optimal λ for a particular distribution

④ Sub-Gaussian

If $X \sim \text{subG}(\sigma^2)$, then we have $E[\exp^{\lambda X}] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right)$ the MGF of $N(0, \sigma^2)$

扩展: If $X \sim N(\mu, \sigma^2)$, then we have $E[\exp^{\lambda X}] = e^{\lambda \mu + \frac{\lambda^2 \sigma^2}{2}}$

$$\text{Proof: } M_X(t) = E[\exp^{\lambda X}] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2} + \lambda x\right) dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-(\mu+\lambda\sigma))^2}{2\sigma^2} + \lambda\mu + \frac{\lambda^2 \sigma^2}{2}\right) dx$$

$$= e^{\lambda\mu + \frac{\lambda^2 \sigma^2}{2}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu-\lambda\sigma)^2}{2\sigma^2}\right) dx = 1$$

$$= e^{\lambda\mu + \frac{\lambda^2 \sigma^2}{2}}$$

⑤ Hoeffding Lemma \rightarrow Hoeffding Bound

Hoeffding Lemma: (Loose Version) $E[\exp^{\lambda(X-E[X])}] \leq \exp\left(\frac{\lambda^2(b-a)}{2}\right)$, $\lambda \in \mathbb{R}$, $X \in \mathbb{Q}$

(Strict Version) $E[\exp^{\lambda(X-E[X])}] \leq \exp\left(\frac{\lambda^2(b-a)}{8}\right)$, $\lambda \in \mathbb{R}$, $X \in [a, b]$

Tight

Proof of the Loose Version:

$$E[\exp^{\lambda(X-E[X])}] = E_X[\exp^{\lambda(X-E_X[X])}] \leq E_X[E_{X'}[\exp^{\lambda(X-X')}]]$$

X is ~~关于~~ 对称的, Rademacher = $\begin{cases} 1, p=\frac{1}{2} \\ -1, p=\frac{1}{2} \end{cases}$, $\Pr(X=x) = \Pr(X=-x)$ $\Rightarrow E_X[E_{X'}[\exp^{\lambda(x-x')}]]$

$$E_{X'}[\exp^{\lambda x}] = \sum_{k=0}^{\infty} \frac{\lambda^k E[x^k]}{k!} \leq E_X[E_{X'}[\exp^{\lambda(x-x')}]]$$

$$= \sum_{k=0,2,4}^{\infty} \frac{\lambda^k E[x^k]}{k!} \quad (\because E[b^k] = 0 \text{ if } k \text{ is odd}) \leq E_X[E_X[\exp^{\lambda^2(a-b)^2}]] \quad (\because a \leq x \leq b)$$

$$= \sum_{k=0,1,2}^{\infty} \frac{\lambda^{2k}}{(2k)!} = \exp\left(\frac{\lambda^2(a-b)^2}{2}\right)$$

$$\leq \sum_{k=0,1,2,\dots}^{\infty} \frac{\lambda^{2k}}{2^k k!} = \sum_{k=0,1,2,\dots}^{\infty} \left(\frac{\lambda^2}{2}\right)^k \frac{1}{k!}$$

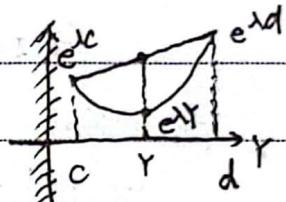
$$= \exp\left(\frac{\lambda^2}{2}\right)$$

Proof of the Tight Version:

Set $Y = X - E[X]$, we have $E[Y] = 0$, $Y \in [a - E[X], b - E[X]] = [c, d]$

What we need to prove is: $E[\exp^{\lambda Y}] \leq \exp\left(\frac{\lambda^2(c-d)^2}{8}\right)$, $\lambda \in \mathbb{R}$, $Y \in [c, d]$, $E[Y] = 0$

D) Proof One (Easy to understand):

$$\begin{aligned} \because e^{\lambda Y} \text{ is convex of } Y \quad \therefore e^{\lambda Y} &\leq \frac{b-y}{b-a} e^{\lambda c} + \frac{y-a}{b-a} e^{\lambda b} \\ E[\exp^{\lambda Y}] &\leq \frac{d-E[Y]}{d-c} e^{\lambda c} + \frac{E[Y]-c}{d-c} e^{\lambda d} \\ &= \frac{d}{d-c} e^{\lambda c} + \frac{-c}{d-c} e^{\lambda d} \\ &= \left(\frac{-c}{d-c}\right) e^{\lambda c} \left(-\frac{d}{c} + e^{\lambda(d-c)}\right) \end{aligned}$$


$\because E[Y] = 0 \quad \therefore c \leq 0, d \geq 0 \text{ and } c, d \text{ don't equal to } 0 \text{ at the same time.} \therefore c < 0, d > 0$

Set $\theta = -\frac{c}{d-c} > 0$,

$$E[\exp^{\lambda Y}] \leq \theta e^{-\lambda \theta(d-c)} \left(\frac{1}{\theta} - 1 + e^{\lambda(d-c)}\right)$$

$$= \underbrace{(1-\theta + \theta e^{\lambda(d-c)})}_{>0} e^{-\lambda \theta(d-c)}$$

$$\because 1-\theta + \theta e^{\lambda(d-c)} = \theta \left(\frac{1}{\theta} - 1 + e^{\lambda(d-c)}\right) = \theta \left(-\frac{a}{b} + e^{\lambda(d-c)}\right) > 0$$

$$\therefore E[\exp^{\lambda Y}] \leq e^{\ln(1-\theta + \theta e^{\lambda(d-c)})} e^{-\lambda \theta(d-c)}$$

~~Set $u = \lambda(d-c)$: $E[e^{\lambda Y}] \leq e^{\ln(1-\theta + \theta e^u) - \theta u}$~~

Define $\varphi: \mathbb{R} \rightarrow \mathbb{R}$, $\varphi(u) = \ln(1-\theta + \theta e^u) - \theta u$, $E[e^{\lambda Y}] \leq e^{\varphi(u)}$

According to Taylor's Theorem, $\exists \xi \in [0, u]$, $\varphi(u) = \varphi(0) + u \varphi'(0) + \frac{1}{2} u^2 \varphi''(\xi)$

$$\varphi(0) = 0, \quad \varphi'(0) = \left(\frac{\theta e^u}{1-\theta+\theta e^u} - \theta\right) \Big|_{u=0} = 0, \quad \varphi''(\xi) = \frac{\theta e^\xi}{1-\theta+\theta e^\xi} \left(1 - \frac{\theta e^\xi}{1-\theta+\theta e^\xi}\right) = \frac{\theta e^\xi}{1-\theta+\theta e^\xi} \leq 0$$

$$\therefore \varphi(u) \leq 0 + 0 + \frac{1}{2} u^2 \times \frac{1}{4} = \frac{1}{8} \lambda^2 (d-c)^2$$

Then we get $E[\exp^{\lambda Y}] \leq \exp\left(\frac{\lambda^2(d-c)^2}{8}\right)$

(2) Proof Two (Hard to understand): $M_Y(\lambda) = E[e^{\lambda Y}]$

Set $\Psi_Y(\lambda) = \ln M_Y(\lambda) = \ln E[e^{\lambda Y}]$,

$$\Psi_Y(0) = \ln E[e^{0 \cdot Y}] = 0, \quad \Psi'_Y(0) = \frac{M'_Y(0)}{M_Y(0)} = \frac{E[Y e^{0 \cdot Y}]}{E[e^{0 \cdot Y}]} = E[Y] = 0$$

For any $s \in [0, \lambda]$,

$$\Psi''_Y(s) = \left(\frac{M'_Y(s)}{M_Y(s)} \right)' = \frac{M''_Y(s)}{M_Y(s)} - \frac{(M'_Y(s))^2}{M_Y(s)^2} = E\left[\frac{Y^2 e^{sy}}{M_Y(s)} \right] - E\left[\frac{Y e^{sy}}{M_Y(s)} \right]^2$$

Define a Z with CDF:

$$F_Z(z) = \int_{-\infty}^z \frac{e^{sy}}{M_Y(s)} dF_Y(y) = \int_c^z \frac{e^{sy}}{E[e^{sy}]} dF_Y(y)$$

F_Z satisfies $F_Z(d) = \int_c^d \frac{e^{sy}}{E[e^{sy}]} dF_Y(y) = 1$, so F_Z is well-defined, and $c \leq z \leq d$

$$E\left[\frac{Y e^{sy}}{M_Y(s)} \right] = \int_c^d y \cdot \frac{e^{sy}}{M_Y(s)} dF_Y(y) = E[Z]$$

$$E\left[\frac{Y^2 e^{sy}}{M_Y(s)} \right] = \int_c^d y^2 \cdot \frac{e^{sy}}{M_Y(s)} dF_Y(y) = E[Z^2]$$

$$\Psi''_Y(s) = E[Z^2] - E[Z]^2 = \text{Var}[Z]$$

$$; c \leq Z \leq d ; |Z - \frac{c+d}{2}| \leq \frac{d-c}{2}$$

$$\therefore \text{Var}[Z] = \text{Var}\left[Z - \frac{c+d}{2}\right] \leq E\left[\left(Z - \frac{c+d}{2}\right)^2\right] = \frac{(d-c)^2}{4}$$

$$\therefore \exists s \in [0, \lambda], \quad \Psi''_Y(s) = \frac{\lambda^2}{2} \text{Var}[Z] \leq \frac{\lambda^2 (d-c)^2}{8}$$

Hoeffding Inequality:

For X_1, X_2, \dots, X_n , if X_i can be bounded. $P(X_i \in [a_i, b_i]) = 1$, Then for

(Independent, r.v.)

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

it satisfies

$$P(\bar{X} - E[\bar{X}] \geq t) \leq \exp\left(-\frac{2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$P(|\bar{X} - E[\bar{X}]| \geq t) \leq 2\exp\left(-\frac{2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad (t > 0)$$

Proof: Using Chernoff Bound, we have

$$P(\bar{X} - E[\bar{X}] \geq t) \leq \min_{\lambda \geq 0} \frac{E[\exp(\lambda(\bar{X} - E[\bar{X}])]}{\exp(\lambda t)}$$

According to Hoeffding Lemma, for $\bar{X} \in [\frac{1}{n} \sum_{i=1}^n a_i, \frac{1}{n} \sum_{i=1}^n b_i]$, we have

$$E[\exp(\lambda(\bar{X} - E[\bar{X}])] \leq \exp\left(\frac{\lambda^2 \frac{1}{n} (\sum_{i=1}^n b_i - \sum_{i=1}^n a_i)^2}{8}\right)$$

so we get

$$P(\bar{X} - E[\bar{X}] \geq t) \leq \min_{\lambda \geq 0} \exp\left(\underbrace{\frac{\lambda^2}{8n^2} (\sum_{i=1}^n b_i - \sum_{i=1}^n a_i)^2}_{g(\lambda)} - \lambda t\right) \stackrel{g'(\lambda)=0}{=} \exp\left(-\frac{2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$\text{Inversely, } P(|\bar{X} - E[\bar{X}]| \geq t) \leq \exp\left(-\frac{2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Finally, we get

$$P(|\bar{X} - E[\bar{X}]| \geq t) \leq 2 \exp\left(-\frac{2t^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Notice: Although Hoeffding Lemma is to prove the Sub-Gaussian property for bounded random variables, we don't need the assumption of bounded Hoeffding Inequality. What we need is just the Sub-Gaussian property, and

$$E[\exp^{\lambda X}] \leq f(x) \quad X \sim \text{sub-G}(\frac{(b-a)^2}{4})$$

③ Bernstein's Inequality (A Sharper bound using the variance of X_i)

let $\{X_i\}$ be independent random variables, where with probability 1. $|X_i - E[X_i]| \leq R$
let $\text{Var}(X_i) \leq \sigma_i^2$. Then, for all $t \geq 0$,

$$P\left(\sum_{i=1}^n (X_i - E[X_i]) \geq t\right) \leq \exp\left(-\frac{t^2}{2\sum \sigma_i^2 + \frac{2}{3}Rt}\right)$$

once this inequality is less used, we omit the proof.

① Azuma - Hoeffding Inequality (Give a bound of martingale difference)

Azuma - Hoeffding aims at solving the dependent variables

Martingale Difference:

In RL, the independent r.v. is very hard to achieve, we usually get the dependent r.v. For dependent r.v., we usually have the martingale difference

Example: Consider two random variables $X_1 \in \{0, 1\}$, $X_2 \in \{-1, 0, 1\}$ with the following distribution:

$X_1=0$	value	probability	$X_1=1$	value	probability
-1	0.50		-1	0	
X_2	0	0	X_2	0	1
1	0.50		1	0	

We have $E[X_2 | X_1=0] = E[X_2 | X_1=1] = E[X_2]$, then we call $\{X_i - E[X_i]\}_{i=1}^n$ is a Martingale difference. Simply explaining, the variables and randomness of the past may change the distribution of the variable I am going to take but can't change the expectation.

Azuma - Hoeffding's Inequality: Suppose F_1, F_2, \dots, F_n are filtrations s.t. $F_{i-1} \subset F_i$, and X_1, X_2, \dots, X_n are random variables s.t. $X_i \in F_i$. If the following conditions hold,

- $E[X_i - E[X_i] | F_{i-1}] = 0$ (Martingale Difference)

- $|X_i - E[X_i]| \leq R_i$,

then we have

$$P\left(\sum_{i=1}^n (X_i - E[X_i]) \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{i=1}^n R_i^2}\right), \forall t \geq 0$$

Proof: What we need is focusing on $E[\exp(\lambda \sum_i (X_i - E[X_i])]$ (Need)

Proof: Set $S_n = \sum_{i=1}^n X_i$, $Z_n = S_n - E[S_n]$, $W_n = e^{\lambda Z_n}$, \mathbb{P}

(Using the proof of Hoeffding lemma)

$$\begin{aligned} E[W_n | \mathcal{F}_{m-1}] &= W_{m-1} E[e^{\lambda(Z_n - Z_{m-1})} | \mathcal{F}_{m-1}] \\ &\leq W_{m-1} \frac{b_n e^{\lambda a_n} - a_n e^{\lambda b_n}}{b_n - a_n} \\ \therefore E[E[W_n | \mathcal{F}_{m-1}]] &\leq E[W_{m-1}] \cdot \frac{b_n e^{\lambda a_n} - a_n e^{\lambda b_n}}{b_n - a_n} \\ &\leq \prod_{i=1}^n \frac{b_i e^{\lambda a_i} - a_i e^{\lambda b_i}}{b_i - a_i} \end{aligned}$$

$$\forall t > 0, P(S_n - E[S_n] \geq t) = P(W_n \geq e^{\lambda t}) \leq \frac{E[W_n]}{e^{\lambda t}}, \forall t > 0$$

$$\leq e^{-\lambda t} \prod_{i=1}^n \frac{b_i e^{\lambda a_i} - a_i e^{\lambda b_i}}{b_i - a_i}$$

$$\leq e^{-\lambda t} \prod_{i=1}^n e^{\lambda^2 (b_i - a_i)^2 / 8}$$

$$\leq \exp\left(-\frac{t^2}{2 \sum_{i=1}^n R_i^2}\right)$$

③ Uniform Concentration

2. (Union Bound) $P(A \cup B) \leq P(A) + P(B)$

2. Suppose $\{X_i\}_{i=1}^n$ are random variables chosen i.i.d., and $\text{supp}(X_i) = \{s \mid P(X_i=s) > 0\} \subseteq S$ where S is a discrete set with $|S| = S$.

Consider two cases shown below.

Case 1: f fixed

$f: S \rightarrow [0,1]$ where $f(x_i)$ are random variables chosen i.i.d. By Hoeffding Ineq.

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i))\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{n}\right)$$

Case 2: f not fixed

For example, $\hat{f} = \arg \min_{f \in F} L(f, \{x_i\}_{i=1}^n)$, \hat{f} depends on the whole dataset

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - E \hat{f}(x_i))\right| \geq t\right) \leq P\left(\exists f \in F, \left|\frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i))\right| \geq t\right)$$

$$= P\left(\bigcup_{f \in F} \left|\frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i))\right| \geq t\right)$$

Case 2a: Finite F

If F is finite, we can use Lemma 1 to get the following:

$$\begin{aligned} P\left(\bigcup_{f \in F} \left|\frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i))\right| \geq t\right) &\leq \sum_{f \in F} P\left(\left|\frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i))\right| \geq t\right) \\ &\leq |F| \exp\left(\frac{-2t^2}{n}\right) \end{aligned}$$

(*) is only true if the following holds with probability $1-\delta$.

$$\sup_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i)) \right| \leq \sqrt{\frac{1}{2n} \log \frac{2|F|}{\delta}}$$

$$= O\left(\sqrt{\frac{1}{n} \log \frac{|F|}{\delta}}\right)$$

case 2b: Infinite F

We can discretize it using ϵ -grid G_ϵ s.t. $G_\epsilon = \{0, \epsilon, 2\epsilon, 3\epsilon, \dots, \lfloor \frac{1}{\epsilon} \rfloor \epsilon\}$

$= G_\epsilon^S = \{f \mid f: S \rightarrow G_\epsilon\}$. We have $\forall f \in F, \exists f \in G_\epsilon$ s.t. $\sup_{x \in S} |f(x) - f_\epsilon(x)| \leq \epsilon$

$$\text{w.p. } 1-\delta, \forall f \in F, \left| \frac{1}{n} \sum_{i=1}^n (f(x_i) - E f(x_i)) \right| \leq \left| \frac{1}{n} \sum_{i=1}^n (f_\epsilon(x_i) - E f_\epsilon(x_i) + (f - f_\epsilon)(x_i) - E[(f - f_\epsilon)(x_i)]) \right|$$

⑨ Bellman Equation

Recall MDP (S, A, r, P, H) where:

environment : P, r

policy : (history independent) $\pi_G : S \rightarrow A / \Delta A$ $\pi_H : H \rightarrow A / \Delta A$

$H : (s_1, a_1, r_1, s_2, a_2, r_2, \dots)$

State-Value Function : $V_h^\pi(s) = E_{\pi_G} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h=s \right]$

Action-Value Function : $Q_h^\pi(s, a) = E_{\pi_G} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h=s, a_h=a \right]$

Bellman Equations :

$$\begin{cases} V_h^\pi(s) = \sum_{a \in A} Q_h^\pi(s, a) \pi_h(a|s) \\ Q_h^\pi(s, a) = r_h(s, a) + E_{s' \sim P_h(s, a)} V_{h+1}^\pi(s') \end{cases}$$

Optimal Policy

Optimality Equation : $V_h^*(s) := \max_{\pi_G} V_h^\pi(s) \quad \forall (s, h)$

Optimality Equation : $Q_h^*(s, a) := \max_{\pi_G} Q_h^\pi(s, a) \quad \forall (h, s, a)$

Optimal policy π^* :

$$\begin{cases} V_h^*(s) := V_h^{\pi^*}(s) \quad \forall (s, h) \\ Q_h^*(s, a) = Q_h^{\pi^*}(s, a) \quad \forall (s, a, h) \end{cases}$$

Theorem 1 (Optimality Existence) : There exists a history independent and deterministic policy π^* that is optimal.

Proof : Let $\pi_H^*(s) = \operatorname{argmax}_a Q^*(s, a)$ be history independent and deterministic

We can now prove by induction :

1) Base case : $H : Q_H^*(s, a) = \pi_H(s, a) = Q_H^{\pi^*}(s, a)$

2) $Q_h^* = Q_h^{\pi^*} \Rightarrow V_h^* = V_h^{\pi^*}$

$$V_h^*(s, a) = \max_{\pi_h} \max_{\pi_{h+1}, \dots, \pi_H} \sum_{a \in A} Q_h^{\pi_h}(s, a) \cdot \pi_h(a|s)$$

$$= \max_{\pi_h} \sum_{a \in A} Q_h^*(s, a) \pi_h(a|s)$$

$$= \max_a Q_h^*(s, a) = Q_h^*(s, \pi_h^*(s)) = V_h^*(s)$$

③ $V_{h+1}^* = V_{h+1}^{\pi_h^*} \Rightarrow Q_h^* = Q_h^{\pi_h^*}$

$$Q_h^*(s, a) = r_h(s, a) + \max_{\pi_h} (P V_{h+1}^{\pi_h})(s, a) = r_h(s, a) + (P V_{h+1}^*)(s, a) = r_h(s, a) + (P V_{h+1}^{\pi_h^*})(s, a) = 0$$

Bellman Optimality Equation: $\begin{cases} V_h^*(s) = \max_a Q_h^*(s, a) \\ Q_h^*(s, a) = r_h(s, a) + (P V_{h+1}^*)(s, a) \end{cases}$

Goal of RL: Find the $\pi_h^*(s) = \arg\max_a Q_h^*(s, a)$

Bellman Equation: $\begin{cases} V^\pi(s) = \sum_{a \in A} Q^\pi(s, a) \pi(a|s) \\ Q^\pi(s, a) = r(s, a) + \gamma (P V^\pi)(s, a) \end{cases}$

② Planning in MDP

Algorithm 1 Policy Evaluation (DP) (Infinite Horizon)

for $h = H, H-1, \dots, 1$ do

$$Q_h^{\pi_h}(s, a) = r_h^{\pi_h}(s, a) + (P_h V_{h+1}^{\pi_h})(s, a)$$

$$V_h^{\pi_h}(s) = \sum_{a \in A} Q_h^{\pi_h}(s, a) \pi_h(a|s)$$

until $\max |V_{h+1}^{\pi_h}(s) - V_h^{\pi_h}(s)| < \theta \times$

Algorithm 2 Policy Iteration

for $t = 1, 2, \dots, H$ do

Run Policy Evaluation to evaluate $Q_h^{\pi^{t-1}}(s, a) \forall s, a, h$

$$\pi_h^t(s) \leftarrow \arg\max_{a \in A} Q_h^{\pi^{t-1}}(s, a)$$

Algorithm 1 Policy Evaluation (Fixed Point)

$$Q_0 \in [0, \beta]^{|A|}$$

$$(r + \gamma + \gamma^2 + \dots + \gamma^\infty) < \frac{1}{1-\gamma} = \beta$$

for $t = 0, 1, \dots$ do

$$Q_{t+1} = T^\pi Q_t = r + \gamma P^\pi Q_t$$

$$\text{Lemma: } \|Q_t - Q^*\|_\infty \leq \gamma^t \|Q_0 - Q^*\|_\infty \leq \gamma^t \beta$$

Algorithm 2 Value Iteration

$$Q_0 \in [0, \beta]^{|A|}$$

for $t = 0, 1, \dots$ do

$$Q_{t+1} = T^* Q_t$$

Algorithm 3 Value Iteration

(Finite Horizon)

(A) Theorem: $T_h^{t^*}$ is optimal $\forall t \leq H$ for $h = H, H-1, \dots, 1$ do

Proof: Using induction.

$$Q_h^*(s, a) = r_h(s, a) + (P_h V_{h+1}^*(s, a))$$

$$\text{Base case: } t=0 \quad Q_H^{t^*}(s, a) = r_H(s, a) = G$$

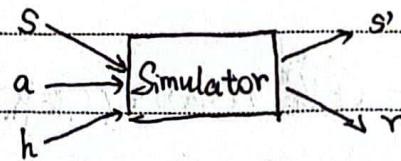
$$V_h^*(s) = \max_{a \in A} Q_h^*(s, a)$$

If t is true, then at $t+1$, $\forall h \leq H-t$,

$$T_h^*(s) = \arg\max_{a \in A} Q_h^*(s, a)$$

$$T_h^{t+1}(s) = \arg\max_a Q_h^{t+1}(s, a) = \arg\max_a Q_h^*(s, a) = T_h^*(s)$$

(II) Generative Models

MDP(S, A, r, P, H) with P, r unknown, but we have access to a sim.

Goal is to find some policy $\tilde{\pi}$ so that $V_i^*(s_i) - \tilde{V}_i^*(s_i) < \epsilon$ for all i , and the sample complexity is how many queries needed to ask the simulator to find the ϵ -optimal $\tilde{\pi}$.

Generative Value Iteration:

Input: n For all $(s, a, h) \in S \times A \times H$ doQuery (s, a, h) and collect n samples $\{s'_1, s'_2, \dots, s'_n\}$

Estimate transition probability by

$$\hat{P}_h(s'|s, a) = \frac{1}{n} \sum_i I[s'_i = s']$$

for all $s' \in S$.For $h = H$ to 1 doCompute $\tilde{Q}_h^*(s, a) = r_h(s, a) + (\hat{P}_h V_{h+1}^*)(s, a)$ Compute $\tilde{V}_h^*(s) = \max_{a \in A} \tilde{Q}_h^*(s, a)$ Output: $\tilde{T}_H^*(s) = \arg\max_a \tilde{Q}_H^*(s, a)$

(Coarse Analysis)

Analysis 1:

Theorem: If we chose $n \geq C \frac{H^4 S}{\epsilon^2} l$, where $l = \log \frac{HS}{P}$, then with probability at least $1 - \epsilon$, the policy resulting from value iteration will be ϵ -optimal.

Lemma 1. For any policy π_G and any $(s, a) \in S \times A$,

$$\underbrace{V_h^\pi(s)}_{\text{true}} - \underbrace{\hat{V}_h^\pi(s)}_{\text{estimated}} = E_{M, \pi_G} \left[\sum_{i=h}^H ((p_i - \hat{p}_i) V_{i+1}^\pi) (s_i, a_i) \mid s_h = s \right]$$

Proof. Note that

$$\begin{aligned} Q_h^\pi(s, a) - \hat{Q}_h^\pi(s, a) &= (PV_h^\pi)(s, a) - (\hat{P}\hat{V}_h^\pi)(s, a) = [(p_h - \hat{p}_h)V_{h+1}^\pi] + (\hat{p}_h)(V_{h+1}^\pi - \hat{V}_{h+1}^\pi) \\ \Rightarrow V_h^\pi(s) - \hat{V}_h^\pi(s) &= E_{\pi_G} [(p_h - \hat{p}_h)V_{h+1}^\pi(s_h, a_h) \mid s_h = s] + E_{M, \pi_G} [V_{h+1}^\pi(s_{h+1}) - \hat{V}_{h+1}^\pi(s_{h+1}) \mid s_h = s] \\ &= E_{M, \pi_G} \left[\sum_{i=h}^H ((p_i - \hat{p}_i) V_{i+1}^\pi) (s_i, a_i) \mid s_h = s \right] \end{aligned}$$

Lemma 2. With probability $\geq 1 - \epsilon$, we have $V(s, a, h) \in S \times A \times [H]$, $\forall V \in [0, H]^S$,

$$|(p_h - \hat{p}_h)V(s, a)| \leq C \cdot H \cdot \frac{SL}{n}, \text{ where } l = \log \frac{HS}{P}$$

Using Lemma 1, 2, Theorem can be get.

(Refined Analysis)

Analysis 2:

Theorem: If we chose $n \geq C \cdot \frac{H^4 l}{\epsilon^2}$ where $l = \log \frac{HS}{P}$, then ...

(12) Q-learning

$$Q_h^t(s,a) = (1-\alpha_t) Q_h^{t-1}(s,a) + \alpha_t (r_h(s,a) + V_{h+1}^{t+1}(s'_t))$$

learning rate

Algorithm Q-learning

Input: learning rate $\{\alpha_t\}_{t=1}^\infty$ Initialize: $Q_h^0(s,a) = 0, \forall (s,a,h) \in S \times A \times [H]$ for $t=1, \dots$ do for $h=1, \dots, H$ do for $(s,a) \in S \times A$ do sample next state s'_t from $P_h(\cdot|s,a)$

Update

$$Q_h^t(s,a) = (1-\alpha_t) Q_h^{t-1}(s,a) + \alpha_t (r_h(s,a) + V_{h+1}^{t+1}(s'_t))$$

 for $s \in S$ do

$$V_h^t(s) = \max_{a \in A} Q_h^t(s,a)$$

Remark:

(1) Model-free learning

2. Incremental update using one noisy sample. Space complexity: $O(HSA)$

From the update rule we get:

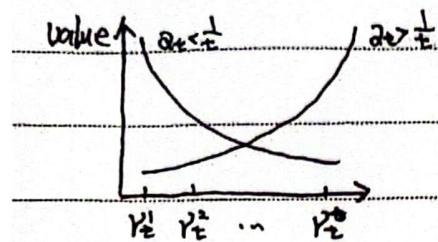
$$Q_h^t(s,a) = \sum_{i=1}^t \gamma_t^i (r_h(s,a) + V_{h+1}^{t+1}(s'_i)) = r_h(s,a) + \sum_{i=1}^t \gamma_t^i V_{h+1}^{t+1}(s'_i)$$

$$\gamma_t^i := \alpha_t \prod_{j=i+1}^t (1-\alpha_j) \text{ and } \sum_{i=1}^t \gamma_t^i = 1$$

(Generative)

Comparing to Value Iteration: $Q_h(s, a) = r_h(s, a) + \sum_{i=1}^n \frac{1}{n} V_{h+1}(s_i')$

1. When $\alpha_t = \frac{1}{t}$ and $\gamma^t = \frac{1}{t}$, Equal
2. When $\alpha_t > \frac{1}{t}$, Q-learning favors later samples.
3. When $\alpha_t < \frac{1}{t}$, Q-learning favors earlier samples.



usually, we choose $\alpha_t > \frac{1}{t}$

⑭ Multi-arm bandit (MAB)

Exploration vs. Exploitation

$$\text{regret}(T) = T \cdot r^* - \sum_{t=1}^T r_{i_t}$$

$$0 \leq \text{regret}(T) \leq T$$

Conversion:

Lemma 1. If algorithm A has regret $c \cdot T^{1-\alpha}$, $\alpha \in (0, 1]$. Then A finds an ϵ -optimal distribution of arm in $(\frac{c}{\epsilon})^{\frac{1}{1-\alpha}}$ samples. Conversely, $c \cdot \epsilon^{-p}$ samples $\geq 2 \cdot c^{\frac{1}{1-p}} \cdot T^{\frac{1}{p}}$

$$\sqrt{T} \text{ Regret} \rightarrow \frac{1}{\epsilon^\alpha} \text{ samples}$$

$$\frac{1}{\epsilon^\alpha} \text{ samples} \rightarrow T^{\frac{2}{\alpha}}$$

Proof: a. $\text{regret}(T) = T \cdot r^* - \sum_{t=1}^T r_{i_t} \leq c \cdot T^{1-\alpha}$

$$\frac{1}{T} \sum_{t=1}^T r_{i_t} \geq r^* - \frac{c}{T^\alpha}$$

$$\frac{c}{T^\alpha} = \epsilon \rightarrow T = (\frac{c}{\epsilon})^{\frac{1}{1-\alpha}}$$

b. $\text{regret}(T) \leq c \cdot \epsilon^{-p} \times 1 + (T - c \cdot \epsilon^{-p}) \times \epsilon$

$$\leq c \cdot \epsilon^{-p} + T \epsilon$$

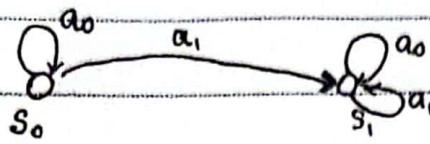
$$= 2 \cdot c^{\frac{1}{1-p}} T^{\frac{p}{1-p}}$$

⑮ Exploration in MDPs

A contextual MAB can be viewed as an MPP without transit

exploration	$\left\{ \begin{array}{l} \text{MABs: } \mathcal{O}(A/\epsilon) \\ \text{MDP: } \mathcal{O}(2^H) \text{ (worst)} \end{array} \right.$
-------------	---

Combinatorial lock:



$$r_H(s_0) = 1, \quad r_H(s_1) = 0$$

$\underbrace{\{a_0, a_1, \dots, a_0\}}_H$ is the \hat{A} optimal, however, we need $\sqrt{2^H}$ to get the po

MPP is more hard for exploration than MAB.

UCB-VI (Upper Confidence Bound Value Iteration)

Algorithm UCB-VI

Initialize: dataset $D := \emptyset$; $Q_h(s, a) := H$ for $h \in [H]$; $V_{H+1}(s) = 0$.

for episode $k=1$ to K do

(Part 1: Value iteration with bonus)

for $h = H$ to 1 do

for $(s, a) \in S \times A$ do

if $(s, a) \in D$ then

Compute $\tilde{P}_h(s'|s, a) := \frac{N_h(s, a, s')}{N_h(s, a)}$ for all $s' \in S$, where $N_h(s, a, s')$ and $N_h(s, a)$

$Q_h(s, a) := \min \{ H, r_h(s, a) + (\tilde{P}_h V_{h+1}(s, a) + b(N_h(s, a))) \};$

for $s \in S$ do

$V_h(s) := \max_{a \in A} Q_h(s, a);$

for $h=1$ to H do

Take action $a_h := \operatorname{argmax}_{a \in A} Q_h(s_h, a)$