Supervised Learning     $(x_1, y_1) \cdots (x_n, y_n)$

$$\min_{\omega} \; L_n(\omega) = \frac{1}{n} \sum_{i=1}^{n} \ell_i \left( g(x_i, \omega), y_i \right)$$

Reinforcement Learning     $(S_t, a_t)$

Multi-Armed Bandit Problem   (MAB-Problem)

Arms        $1, 2, \cdots, n$

Rewards     $D_1, D_2, \cdots, D_n \in [0, 1]$

a-priori unknown

$\mu_i = \mathbb{E} D_i \quad i = 1, \cdots, n$

Bandit

1  2  3  4  $\cdots$  n

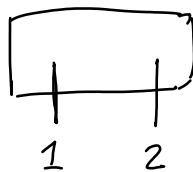Target          Find $i^* = \arg\max\limits_{i = 1 \cdots n} \{M_i\}$   ?

Arm sequence played $= i_1, i_2, \ldots, i_t, \cdots$

Reward Sequence $=$ $\quad r_1, r_2, \cdots, r_t, \cdots$

$$r_t = \text{an i.i.d. copy of } D_{i_t}$$

Total Reward   $\sum\limits_{t=1}^{T} r_t$

Suppose    $n = 2$



$$D_1 \qquad D_2$$
$$\uparrow \qquad \uparrow$$
$$\mu_1 = \mathbb{E} D_1 \qquad \mu_2 = \mathbb{E} D_2$$

Naive Algorithm

$R_i^{(k)}$ are i.i.d. rewards each time $^{(k)}$ when the $i$-th arm is played

$R_i^{(k)}$ is a copy of $D_i$

Step 1    play arm 1    $T/2$ times

play arm 2    $T/2$ times

$$\widehat{\mu_1} = \frac{R_1^{(1)} + R_1^{(2)} + R_1^{(3)} + \ldots + R_1^{(T/2)}}{T/2}$$

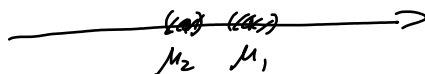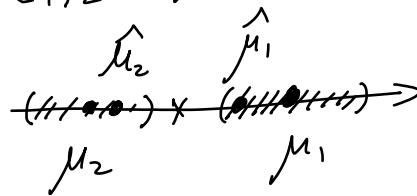$$\xrightarrow[L.L.N]{T/2} \quad \mathbb{E} R_1^{(1)} = \mu_1 \qquad \boxed{|\widehat{\mu_i} - \mu_i|}$$

similarly    $\widehat{\mu_2} = \dfrac{R_2^{(T/2+1)} + \ldots + R_2^{(T)}}{T/2} \xrightarrow{L.L.N.} \mathbb{E} R_2^{(1)} = \mu_2$

$\underline{Step 2}$     $\hat{i}_T = \arg\max_{i=1,2} \{\hat{\mu_i}\}$

Assume   $\mu_1 > \mu_2$

Arm 1 is best

$\hat{\mu_2}$   $\hat{\mu_1}$

$\mu_2$      $\mu_1$

$\mu_2$  $\mu_1$

intervals $\rightarrow 0$

as $T \rightarrow \infty$

"concentration"

$\Delta = |\mu_1 - \mu_2|$

**Theorem**     Naive — Thinking Algorithm Outputs better arm

$$i^* = \arg\max_{i = 1, 2} \{\mu_i\}$$

with probability $\geq 1 - 4\exp\left(-\dfrac{\Delta^2 T}{4}\right)$

(Hoeffding's Ineq.) If $X_1, X_2, \ldots X_n$ are i.i.d. r.v.'s

$X_i \in [a_i, b_i]$ for $\forall i \in \{1 \ldots n\}$

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

$$\mathbb{P}\left[\,|\mu_i - \hat{\mu_i}| \leq \frac{\Delta}{2}\,\right] = 1 - \mathbb{P}\left[\,|\mu_i - \hat{\mu_i}| > \frac{\Delta}{2}\,\right]$$

$$\mathbb{P}(A^c)$$
$$\mathbb{P}(B^c) \leq 2\exp\left(\frac{-\Delta^2 T}{4}\right)$$

$$\geq 1 - 2\exp\left(\frac{-\Delta^2 T}{4}\right)$$

$$A = \left\{\,|\mu_1 - \hat{\mu_1}| \leq \frac{\Delta}{2}\,\right\} \qquad B = \left\{\,|\mu_2 - \hat{\mu_2}| \leq \frac{\Delta}{2}\,\right\} \qquad \mathbb{P}(A \cap B)$$

$$\mathbb{P}(A \cap B) = 1 - \mathbb{P}\left((A \cap B)^c\right) = 1 - \mathbb{P}\left(A^c \cup B^c\right)$$

$$\geq 1 - \mathbb{P}(A^c) - \mathbb{P}(B^c)$$

$$\mathbb{P}(A^c \cup B^c) \leq \mathbb{P}(A^c) + \mathbb{P}(B^c)$$

$$\geq 1 - 4\exp\left(-\frac{\Delta^2 T}{4}\right)$$

maximizing the expected overall rewards $\quad \mathbb{E}\left(\sum_{t=1}^{T} r_t\right)$

Regret $\qquad R_T = T \max\{\mu_1 \cdots \mu_n\} - \mathbb{E}\left(\sum_{t=1}^{T} r_t\right)$

$$= \mathbb{E}\left[\sum_{t=1}^{T}\left(\max\{\mu_1 \cdots \mu_n\} - r_t\right)\right]$$

Problem.        minimize regret $\cdots$

Example    For Naive-Thinking Algorithm   $\mu_1 > \mu_2$

$$R_T = \frac{T}{2}(\mu_1 - \mu_2) = \frac{T}{2}\Delta \sim \mathcal{O}(T) \quad \left(\begin{array}{l}\mathcal{O}(T \ln T) \\ \mathcal{O}(\ln T) \\ \cdots \end{array}\right)$$

Upper - Confidence - Bound  (UCB) method

How to understand Regret?

Assume there are $n$-arms    $\mu_1 > \mu_2 \geq \cdots \geq \mu_n$

$$\Delta_i = \mu_1 - \mu_i > 0$$

$$\sum_{t=1}^{T} r_t = \sum_{i=1}^{n} \sum_{k=1}^{T_i(T)} R_i^{(k)} = \sum_{i=1}^{n} T_i \underbrace{\frac{\sum_{k=1}^{T_i(T)} R_i^{(k)}}{T_i}}_{= \widehat{\mu}_{i,T}} = \sum_{i=1}^{n} T_i \widehat{\mu}_{i,T}$$

$T_i = T_i(t) = \#$ of times that the $i$-th arm is played up to time $t$

$$\text{Regret} = R_T = \mathbb{E}\left[ \sum_{t=1}^{T} (\mu_1 - r_t) \right] \qquad T = T_1 + \cdots + T_n$$

$$= T_1(T) + \cdots + T_n(T)$$

$$= T \mu_1 - \mathbb{E}\left( \sum_{t=1}^{T} r_t \right)$$

$$= T \mu_1 - \mathbb{E}\left( \sum_{i=1}^{n} T_i \,\widehat{\mu_{i},T} \right)$$

$$= (T_1 + \cdots + T_n) \mu_1 - \mathbb{E}\left( \sum_{i=1}^{n} T_i \,\widehat{\mu_{i},T} \right)$$

$$= \mathbb{E}\left( \sum_{i=1}^{n} T_i \,( \mu_1 - \widehat{\mu_{i},T} ) \right) \qquad \frac{\mu_1 - \widehat{\mu_{i},T}}{\boxed{T_i}} \,?$$

$$\text{Regret} = R_T = \mathbb{E}\left[ \sum_{i=1}^{n} T_i \left( \mu_1 - \widehat{\mu}_{i,T} \right) \right] = \mathbb{E}\left[ \sum_{i=1}^{n} T_i (\mu_1 - \mu_i) \right]$$
$$+ \mathbb{E}\left[ \sum_{i=1}^{n} T_i (\mu_i - \widehat{\mu}_{i,T}) \right]$$

$$\mu_1 > \mu_2 \geqslant \cdots \geqslant \mu_n \qquad \mathbb{E}\,\widehat{\mu}_{i,T} = \mu_i \quad \in [0,1]$$

$$\widehat{\mu}_{i,T} = \frac{1}{T_i} \sum_{k=1}^{T_i} R_i^{(k)}$$

$$\mathbb{E}\left( \mu_1 - \widehat{\mu}_{i,T} \right) = \mu_1 - \mu_2 \qquad\qquad \mu_1 > \mu_i$$

To make $R_T$ small  want to play arm 1 as many as you can

But this make $T_1$ large !

Hoeffding's Inequality

$$\mathbb{P}\left[\left|\mu_i - \widehat{\mu}_{i,T}\right| \leq \sqrt{\frac{\ln(2Tn)}{T_i}}\right]$$

$$= 1 - \mathbb{P}\left[\left|\mu_i - \widehat{\mu}_{i,T}\right| > \sqrt{\frac{\ln(2Tn)}{T_i}}\right]$$

$$= 1 - \mathbb{P}\left[\left|\mu_i - \frac{\sum_{k=1}^{T_i} R_i^{(k)}}{T_i(t)}\right| > \sqrt{\frac{\ln(2Tn)}{T_i}}\right]$$

$$\geq 1 - 2\mathbb{P}\left[\frac{\sum_{k=1}^{T_i} R_i^{(k)}}{T_i} - \mu_i > \sqrt{\frac{\ln(2Tn)}{T_i}}\right]$$

$$\mathbb{P}\left(|a-b| > c\right)$$

$$= \mathbb{P}\left(a-b > c \text{ or } b-a > c\right)$$

$$\leq \mathbb{P}(a-b>c) + \mathbb{P}(b-a>c)$$

$$X_k = \frac{R_i^{(k)}}{T_i}$$

$$\in \left[0, \frac{1}{T_i}\right]$$

$$1 - 2\exp\left(-\frac{2 \cdot \frac{\ln(2Tn)}{T_i}}{\sum_{k=1}^{T_i} \frac{1}{T_i^2}}\right) \quad \frac{1}{T_i}$$

$$\mathbb{P}\left[\, |\mu_i - \widehat{\mu_{i,T}}| \le \sqrt{\frac{\ln(2Tn)}{T_i}} \,\right] \ge 1 - 2 \cdot \left(\frac{1}{2Tn}\right)^2 = 1 - \frac{1}{2T^2n^2}$$

"concentration"

i.e.   $\widehat{\mu_{i,T}} \in \left[\, \mu_i - \sqrt{\frac{\ln(2Tn)}{T_i}}, \; \mu_i + \sqrt{\frac{\ln(2Tn)}{T_i}} \,\right]$

confidence interval

with probability $\ge 1 - \frac{1}{2T^2n^2}$

$$\underline{UCB - Arm} \qquad \hat{i}_t = \arg\max_i \left\{ \underline{\hat{\mu}_{i,t-1}} + \sqrt{\frac{\ln(2Tn)}{T_i(t-1)}} \right\}$$

$$Set \quad \overline{E}_i = \left\{ |\mu_i - \hat{\mu}_{i,T}| \leq \sqrt{\frac{\ln(2Tn)}{T_i}} \right\}$$

$$E = \bigcap_{i=1}^{n} \overline{E}_i \qquad \mathbb{P}(E) \geq 1 - \sum_{i=1}^{n} \mathbb{P}(\overline{E}_i^C)$$

$$\geq 1 - \sum_{i=1}^{n} \frac{1}{2T^2 n^2} = 1 - \frac{1}{2T^2 n}$$

$$\mathbb{P}(E^C) \leq \frac{1}{2T^2 n}$$

Let $t_i$ be the last time that the arm $i$ is played

on $\overline{E}$ $\qquad \mu_i + 2\sqrt{\dfrac{\ln(2Tn)}{T_i(t_i-1)}} \overset{\text{on } \overline{E}}{\geqslant} \widehat{\mu}_{i,\,t_i-1} + \sqrt{\dfrac{\ln(2Tn)}{T_i(t_i-1)}}$

$\overset{UCB}{\geqslant} \widehat{\mu}_{1,\,t_i-1} + \sqrt{\dfrac{\ln(2Tn)}{T_1(t_i-1)}}$

$\overset{\text{on } \overline{E}}{\geqslant} \mu_1$

$$\mu_i + 2\sqrt{\frac{\ln(2T_n)}{T_i(t_i-1)}} \geq \mu_1 > \mu_{\bar{i}} \quad \text{on } \bar{E}$$

$$\Leftrightarrow \quad \frac{\ln(2T_n)}{T_i(t_i-1)} \geq \left(\frac{\mu_1 - \mu_i}{2}\right)^2$$

$$\Leftrightarrow \quad T_i \leq \frac{4\ln(2T_n)}{(\mu_1-\mu_i)^2} \sim \frac{\ln(T_n)}{\Delta_i^2}$$

$$\mu_1 - \mu_i = \Delta_i$$

$$R_T^{UCB} = \mathbb{E}\left[\sum_{i=1}^{n} T_i \left(\mu_1 - \widehat{\mu_{i,T}}\right)\right]$$

$$\mu_i \in [0,1]$$
$$R_i^{(k)} \in [0,1]$$

$$\leq \mathbb{E}\left[\sum_{i=1}^{n} T_i \left(\mu_1 - \widehat{\mu_{i,T}}\right) \Big| E\right] + \mathbb{P}(E^C) \cdot T$$

$$\boxed{R_T^{UCB} \lesssim \mathcal{O}(\sqrt{nT\ln T})}$$

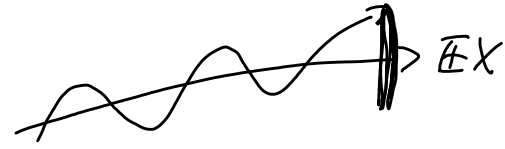$$\mathbb{E}\,\widehat{\mu_{i,T}} = \mu_i$$

$$\Delta_i = \mu_1 - \mu_i$$

$$\leq \mathbb{E}\left[\sum_{i=2}^{n} T_i \Delta_i \Big| E\right] + \mathbb{P}(E^C) \cdot T$$

$$\lesssim \mathbb{E}\left[\sum_{i=2}^{n} \sqrt{T_i} \cdot \sqrt{\frac{\ln(nT)}{\Delta_i^2}} \Delta_i \Big| E\right] + \underbrace{\mathbb{P}(E^C) \cdot T}_{\frac{1}{2T^3 n} \cdot T \sim \mathcal{O}(\frac{1}{T})}$$
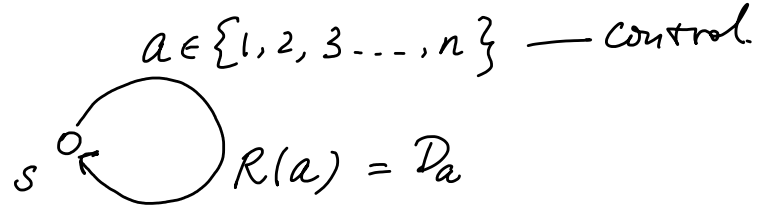
$$= \mathbb{E}\left[\sum_{i=2}^{n} \sqrt{T_i \ln(nT)} \Big| E\right] \qquad T$$

$$\leq \mathbb{E}\left[(n-1)\sqrt{\ln(nT)\left(\sum_{i=1}^{n}\frac{T_i}{n-1}\right)} \Big| E\right] \sim \sqrt{(n-1)\ln(nT)T} \sim \sqrt{nT\ln T}$$

" Exploration — Exploitation "

"online learning"

EX

Bandit

$a \in \{1, 2, 3 \ldots, n\}$ — Control

$s$     $R(a) = D_a$

Markov Decision Processes
        (MDP)

$$\underline{MDP} = MDP(\mathcal{S}, \mathcal{A}, H, \underline{\mathbb{P}}, r)$$

$\mathcal{S}$ —— set of states  $|\mathcal{S}| = S$

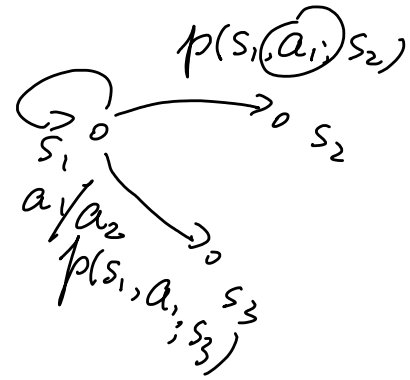$\mathcal{A}$ —— set of actions  $|\mathcal{A}| = A$

$H$ —— number of steps in each episode  "finite horizon"

$\underline{\mathbb{P}}$ —— transition matrix  "under control by $\mathcal{A}$"

$\mathbb{P}_h(\cdot | x, a)$ is the transition probabilities at state $x$ when action $a$ is taken at step $h \in [H]$

$r_h : \mathcal{S} \times \mathcal{A} \to [0,1]$
is a deterministic reward at step $h \in [H]$

$p(s_1, \boxed{a_1}, s_2)$

$s_1 \circ \quad \to^{\circ} s_2$

$a_1 / a_2$

$p(s_1, a_1, s_3 ; s_3)$

Dynamics of our MDP

In each episode of MDP   initial state $x_1$ is picked arbitrarily

at step $h \in [H]$ of episode

agent observes $x_h \in \mathcal{S}$

picks action $a_h \in \mathcal{A}$

receive reward $r_h(x_h, a_h)$

then transit to next state $x_{h+1}$

drawn from $\mathbb{P}_h(\cdot | x_h, a_h)$

episode ends when $x_{H+1}$ is reached

Policy    $\pi = \left\{ \pi_h : \mathcal{S} \to \mathcal{A} \right\}_{h \in [H]}$

"at step $h$ take policy $\pi_h$ take $a_h = \pi_h(x_h) \in \mathcal{A}$.

Value function    $V_h^{\pi} : \mathcal{S} \to \mathbb{R}$   is the value function at step $h$ under policy $\pi$

$$V_h^{\pi}(x) = \mathbb{E}\left[ \sum_{h'=h}^{H} r_{h'}(x_{h'}, \pi_{h'}(x_{h'})) \,\middle|\, x_h = x \right]$$

Optimal Policy $\pi^*$    $V_h^{\pi^*}(x) \equiv V_h^*(x) = \sup_{\pi} V_h^{\pi}(x)$ for all $x \in \mathcal{S}$ and $h \in [H]$

Regret        Agent is to play this MDP of $K$ episodes

$$k = 1 \cdots K$$

suppose we pick a starting state $x_1^k$ for each episode $k$

$$Regret(K) = \sum_{k=1}^{K} \left[ V_1^*(x_1^k) - V_1^{\pi_k}(x_1^k) \right]$$

$$\pi_k = \left\{ \pi_{k,h} : \mathcal{S} \to \mathcal{A} \right\}$$

Key        Q-value function at step $h$

$$Q_h^\pi: \quad \mathcal{S} \times \mathcal{A} \longrightarrow \mathbb{R}$$

$$Q_h^\pi(x, a) = r_h(x, a) + \mathbb{E}\left[\sum_{h'=h+1}^{H} r_{h'}(x_{h'}, \pi_{h'}(x_{h'})) \,\middle|\, \begin{matrix} x_h = x \\ a_h = a \end{matrix}\right]$$

$$\pi_h^*(x) = \arg\max_a Q_h^\pi(x, a)$$

$$\pi_h^* \quad .$$

Bellmann's equations

$$\begin{cases} V_h^\pi (x) = Q_h^\pi (x, \pi_h(x)) \\ \\ Q_h^\pi (x, a) = r_h(x, a) + \overline{\mathbb{E}}_{x' \sim \mathbb{P}_h(\cdot | x, a)} V_{h+1}^\pi (x') \\ \\ V_{H+1}^\pi (x) = 0 \qquad \forall x \in \mathcal{S} \end{cases}$$

Dynamical Programming

$$Q_h^\pi (x, a) : \underset{\{1 \cdots H\}}{[H]} \times \mathcal{S} \times \mathcal{A} \longrightarrow \mathbb{R}$$

$\pi^* \longrightarrow$ optimal policy

Bellmann's Optimality Equations

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

$$Q_h^*(x, a) = r_h(x, a) + \overline{\mathbb{E}}_{x' \sim \mathbb{P}_h(\cdot \mid x, a)} V_{h+1}^*(x')$$

$$V_{H+1}^*(x) = 0 \qquad \forall \, x \in \mathcal{S}$$

$$\pi^* \longrightarrow Q_h^{\pi^*}(x, a)$$

$\underline{Ex}$
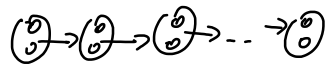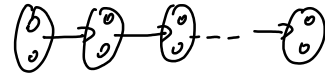
$S_1, S_2, \cdots, S_n, \cdots$

$\pi_1 = \dfrac{p_{21}}{p_{12}+p_{21}}$    $\pi_2 = \dfrac{p_{12}}{p_{12}+p_{21}}$    Monte-Carlo

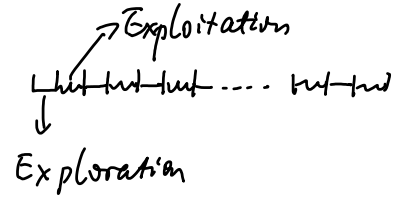$\pi_i \approx \dfrac{\# \text{ of visits to state } i}{n}$    $i = 1, 2$

Monte-Carlo Dynamical Programing

$$Q_h^*(x, a) \approx r_h(x, a) + \underbrace{\dfrac{V_{h+1}^*(x_{h+1}^{(1)}) + V_{h+1}^*(x_{h+1}^{(2)}) \cdots + V_{h+1}^*(x_{h+1}^{(N)})}{N}}_{SS}$$

Exploration  Exploitation

$$\mathbb{E}_{x' \sim \mathbb{P}_h(\cdot | x, a)} V_{h+1}^*(x').$$

## Q - learning.

$$\frac{r_1 + r_2 + \cdots + r_n}{n} = Q_n \approx \overline{\underline{E\ V}}$$

$$\frac{r_1 + r_2 + \cdots + r_{n+1}}{n+1} = Q_{n+1}$$

$$Q_{n+1} = \frac{r_{n+1}}{n+1} + \frac{r_1 + \cdots + r_n}{n+1} = \frac{r_{n+1}}{n+1} + \frac{n}{n+1}\left(\frac{r_1 + \cdots + r_n}{n}\right)$$

$$Q_{n+1} = \frac{n}{n+1}Q_n + \frac{r_{n+1}}{n+1} = Q_n + \underbrace{\left(\frac{1}{n+1}\right)}_{\alpha_n}\underbrace{(r_{n+1} - Q_n)}_{\text{"Incremental"}} \quad TD$$

Exploitation

Exploration

## $Q$-learning iteration

for episode $k = 1 \cdots K$    do

$\quad$ receive $x_1$
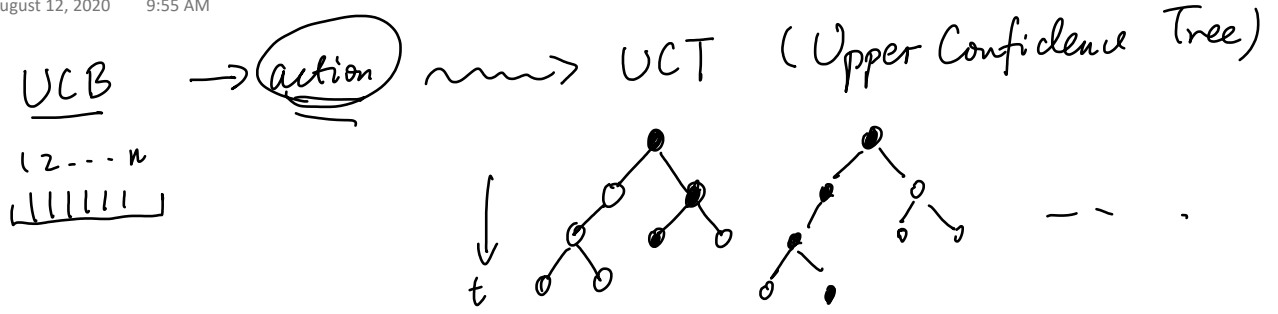
$\quad$ for step $h = 1 \cdots H$    do

$\qquad a_h \leftarrow \text{argmax}_{a'} Q_h(x_h, a')$    observe $x_{h+1}$

$\qquad t = N_h(x_h, a_h) \leftarrow N_h(x_h, a_h) + 1$

$\qquad Q_h(x_h, a_h) \leftarrow Q_h(x_h, a_h) + \alpha_t \left[ r_t(x_h, a_h) + V_{h+1}(x_{h+1}) \phantom{\Big]} \right. $
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad - Q_h(x_h, a_h) $

$\qquad V_h(x_h) \leftarrow \min \left\{ H , \max_{a' \in \mathcal{A}} Q_h(x_h, a') \right\}$

UCB $\longrightarrow$ (action) $\rightsquigarrow$ UCT (Upper Confidence Tree)

$1\,2\cdots n$



$t$

Jin C. Zhu. Z.A. Bubeck.S & Jordan. M.   Is Q-learning provably efficient?

NeuIPS 2018

$$Q_h(x_h, a_h) \leftarrow Q_h(x_h, a_h) + \alpha_t \left[ r_h(x_h, a_h) + V_{h+1}(x_{h+1}) - Q_h(x_h, a_h) \right]$$

$\sqrt{\frac{1}{t}} \quad +b_t$
$\phantom{xx} s$