# UCB-VI Algorithm

Recap :    So far,    MDP w/ generative model (simulator)

   Today:   exploration in MDP

Setup - $M = (S, A, H, \{P_h\}_{h \in [H]}, \{r_h\}_{h \in [H]})$

  - no assumption on generative model
  - interaction protocol:

for each game (episode):
   - start of the episode : $s_1 \sim d_1(s)$
   - for $h = 1, 2, .., H$ :
      · take $a_h \in A$
      · observe $s_{h+1} \sim P_h(\cdot | s_h, a_h)$ and reward $v_h = r_h(s_h, a_h)$

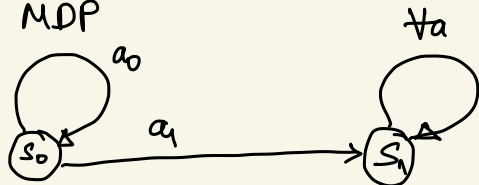regret minimization : Find sequence of policies $\{\pi_k\}_{k \in [K]}$

to minimize $\text{regret}(K) = K \cdot V_1^*(s_1) - \sum_{k=1}^{K} V_1^{\pi_k}(s_1)$

Review    Story from MAB

| Alg | Regret |
|---|---|
| explore-the-commit | $\tilde{O}\left(A^{1/3} T^{2/3}\right)$ |
| $\varepsilon$-greedy | $\tilde{O}\left(A^{1/3} T^{2/3}\right)$ |
| UCB | $\tilde{O}\left(\sqrt{AT}\right)$ |
| lower bound | $\Omega\left(\sqrt{AT}\right)$ |

**Lemma**  Random exploration requires $\Omega(2^H)$ samples to find an optimal policy

**Proof:**  "Combinatorial lock" MDP



zero rewards everywhere except at $(H, S_0)$ where $f_H(S_0, a) = 1$  $\forall a$

Optimal policy: $\qquad \pi_h^*(S_0) = a_0$ $\qquad,\qquad V_1^{\pi^*}(S_0) = 1$

$\qquad\qquad\qquad\qquad \pi_h^*(S_1) = a_0 \text{ or } a_1 \qquad V_1^{\pi^*}(S_1) = 0$

To discover $\pi^*$, action sequence must be $(a_0, \ldots, a_0)$

$$Pr_{\mathcal{U}M_S}\Big((a_0, \ldots, a_0)\Big) = \frac{1}{2^H}$$

$\Rightarrow$ need $\Omega(2^H)$ episodes to discover $(a_0, \ldots, a_0)$ for once

# UCB - VF

Let $b(N) = c \sqrt{\frac{H^2 L}{N}}$ where $L := \log\left(\frac{SAHK}{\delta}\right)$

- initialize $D = \emptyset$, $Q_h(s,a) = 0$ $\forall h \in [H]$ $Q_{H+1}(s) = 0$

- for $k = 1, 2, \ldots, K$:   (estimation phase)

  • $\widehat{P}_h(s'|s,a) = \dfrac{N_h(s,a,s')}{N_h(s,a)}$

  where $N_h(s,a,s') = \left| \{ (h,s,a,s') \in D \} \right|$

  $N_h(s,a) = \left| \{ (h,s,a) : (h,s,a,s') \in D \} \right|$

  • $Q_h(s,a) = \left[ r_h(s,a) + \left( \widehat{P}_h V_{h+1} \right)(s,a) + \underbrace{b(N_h(s,a))}_{\text{bonus}} \right]_{[0,H]}$

  • $V_h(s) = \max_a Q_h(s,a)$

- (Excecution phase)

  For $h = 1, \ldots, H$:

  Take $a_h = \underset{a}{\text{argmax}}\ Q_h(s,a)$

  Observe $s_{h+1}$, $r_h$

  $D = D \cup \{ (h, s_h, a_h, s_{h+1}) \}$

**Theorem**  w.p. al $1-\delta$, the regret of UCB-VI is:

$$\text{regret}(K) \leq c \cdot \left( H^2 \sqrt{SAK\iota} + H^3 S^2 A \iota^2 \right)$$

**Notations**  Add superscript $k$ to all quantities

- $D^k$:  $D$ up to $k$-th episode

- $N_h^k(s,a,s')$:  number of $(h,s,a,s')$ in $D^k$

- $\hat{P}_h^k$:  empirical distribution

- $(s_h^k, a_h^k, s_{h+1}^k)$:  tuple played at $k$-th episode

**Lemma** (optimism)   wpal $1-\delta$:

$$Q_h^k(s,a) \geq Q_h^*(s,a), \qquad V_h^k(s) \geq V_h^*(s) \qquad \forall (k,h,s,a)$$

**Proof** (by induction)

- when $h = H+1 \longrightarrow$ trivial
- Assume by induction that it holds for some $h+1$.
- $Q_h^k(s,a) - Q_h^*(s,a) = \left(\hat{P}_h^k V_{h+1}^k\right)(s,a) + b(N_h^k(s,a)) - \left(P_h V_{h+1}^*\right)(s,a)$

$$= \hat{P}_h^k \left(V_{h+1}^k - V_{h+1}^*\right)(s,a)$$

$$+ \underbrace{\left(\hat{P}_h^k - P_h\right) V_{h+1}^*(s,a) + b(N_h^k(s,a))}_{\geq 0 \quad \text{by Hoeffding's ineq}}$$

- $V_h^k(s) = \max\limits_a Q_h^k(s,a) \geq \max\limits_a Q_h^*(s,a) = V_h^*(s)$

$$\cdot \text{ regret }(K) = \sum_{k=1}^{K} \left( V_1^{\pi^{**}}(s_1) - V_1^{\pi_k}(s_1) \right) \leq \sum_{k=1}^{K} \left( V_1^{k}(s_1) - V_1^{\pi_k}(s_1) \right)$$

Optimism remove the unknown $\pi^{**}$ from our bound

$\longrightarrow$ make our job easier

$$V_h^k(s_1^k) - V_h^{\pi_k}(s_1^k) =$$

$$(Q_h^k - Q_h^{\pi_k})(s_h^k, a_h^k) \leq (\hat{P}_h^k V_{h+1}^k - P_h V_{h+1}^{\pi_k})(s_h^k, a_h^k) + b_h^k$$

$$= (\hat{P}_h^k - P_h) V_{h+1}^k (s_h^k, a_h^k) + P_h(V_{h+1}^k - V_{h+1}^{\pi_k})(s_h^k, a_h^k) + b_h^k \qquad b(N_h^k(s_h^k, a_h^k))$$

$$= \underbrace{(\hat{P}_h^k - P_h) V_{h+1}^\star (s_h^k, a_h^k)}_{\leq b_h^k} + (\hat{P}_h^k - P_h)(V_{h+1}^k - V_{h+1}^\star)(s_h^k, a_h^k)$$

$$\qquad\qquad\qquad + P_h(V_{h+1}^k - V_{h+1}^{\pi_k})(s_h^k, a_h^k) + b_h'^k$$

$$\leq \underbrace{(\hat{P}_h^k - P_h)(V_{h+1}^k - V_{h+1}^\star)(s_h^k, a_h^k)}_{I} + \underbrace{P_h(V_{h+1}^k - V_{h+1}^{\pi_k})(s_h^k, a_h^k)}_{J} + 2b_h^k$$

Note: $(\hat{P}_h^k - P_h) V(s_h^k, a_h^k) \leq b_n^k$  but  $V_{h+1}^k$  is data-dependent!

$$U(s') = \left( V_{h+1}^k - V_{h+1}^* \right)(s') \quad \forall s' \in S$$

$$I = \sum_{s' \in S} \left( \hat{P}_h^k(s' \mid s_h^k, a_h^k) - P_h(s' \mid s_h^k, a_h^k) \right) U(s')$$

$$\leq c \sum_{s' \in S} \left[ \sqrt{\frac{P_h(s' \mid s_h^k, a_h^k) \, \iota}{N_h^k(s_h^k, a_h^k)}} + \frac{\iota}{N_h^k(s_h^k, a_h^k)} \right] \cdot U(s') \quad \text{(Bernstein's)}$$

$$\leq c \sum_{s' \in S} \left[ \frac{P_h(s' \mid s_h^k, a_h^k)}{cH} + \frac{cH\iota}{N_h^k(s_h^k, a_h^k)} \right] U(s') \quad \text{(AM-GM)}$$

$$= \boxed{\frac{1}{H}} P_h \left( V_{h+1}^k - V_{h+1}^* \right)(s_h^k, a_h^k) + c \frac{S H^2 \iota}{\underbrace{N_h^k(s_h^k, a_h^k)}_{S_h^k}}$$

$$V_h^k(s_n^k) - V_h^{\pi_k}(s_n^k) \leq I + J + 2b_h^k$$

$$\underbrace{\phantom{V_h^k(s_n^k) - V_h^{\pi_k}(s_n^k)}}_{\Delta_h^k}$$

$$\leq \frac{1}{H} P_h (V_{h+1}^k - V_{h+1}^*)(s_n^k, a_h^k) + S_n^k$$

$$P_h(V_{h+1}^k - V_{h+1}^{\pi_k})(s_h^k, a_h^k) + 2b_h^k$$

$$\vdots$$

$$\leq \left(1 + \frac{1}{H}\right) \underbrace{P_h (V_{h+1}^k - V_{h+1}^{\pi_k})(s_h^k, a_n^k)}_{\xi_h^k + (V_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1}^k)} + 2b_h^k + S_h^k$$

where

$$\xi_h^k = P_h(V_h^k - V_{h+1}^{\pi_k})(s_n^k, a_h^k) - (V_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1}^k)$$

(martingale).

$$\Delta_h^k \leq \left(1+\frac{1}{H}\right)\left(\Delta_{h+1}^k + \xi_h^k\right) + 2b_h^k + S_h^k$$

$$= \left(1+\frac{1}{H}\right)\Delta_{h+1}^k + \left(1+\frac{1}{H}\right)\xi_h^k + S_n^k + 2b_n^k$$

$$\Delta_1^k \leq \underbrace{\left(1+\frac{1}{H}\right)^H}_{e}\Delta_H^k + \underbrace{\left(1+\frac{1}{H}\right)^H}_{\leq e}\sum_{h=1}^{H}\left(\xi_h^k + S_h^k + b_h^k\right)$$

$$\Rightarrow \text{regret}(K) = \sum_{k=1}^{K}\Delta_1^k \leq c.\sum_{k=1}^{K}\sum_{h=1}^{H}\left(\xi_n^k + S_n^k + b_n^k\right)$$

$$\bullet \quad \sum_k \sum_n b_n^K = {}^{C_i} H\sqrt{i} \sum_k \sum_n \frac{1}{\sqrt{N_h^{K}(s_{h,}^k a_h^k)}}$$

$$= c \cdot H\sqrt{i} \sum_h \sum_{(s,a)} \sum_{i=1}^{N_h^{K}(s,a)} \frac{1}{\sqrt{i}}$$

$$\leq c H\sqrt{i} \sum_{(s,a,h)} \sqrt{N_h^{K}(s,a)}$$

$$= c H \sqrt{i} \sqrt{SAH} \sqrt{KH} \qquad = H^2 \sqrt{SAi c}$$

$$\bullet \quad \sum_k \sum_n \delta_n^K \leq c \cdot SH^2 i \sum_{k,h} \frac{1}{N_h^{K}(s_h^k, a_h^k)}$$

$$\leq c \cdot SH^2 i \sum_{(h,s,a)} \sum_{i=1}^{N_h^{K}(s,a)} \frac{1}{i}$$

$$\leq c \cdot SH^2 i \sum_{(h,s,a)} \log N_h^{K}(s,a) \leq c \cdot SH^2 i \log (KH)$$

- $\displaystyle\sum_k \sum_n \xi_n^k \quad \leqslant H^2 \cdot \sqrt{K}$