Solution by FILL fullname command at the beginning of latex document

**Instructions**

- The deadline is **January 10, 2021. 23h00**

- By doing this homework you agree to the *late day policy, collaboration and misconduct rules* reported on Piazza.

- **Mysterious or unsupported answers will not receive full credit**. A correct answer, unsupported by calculations, explanation, or algebraic work will receive no credit; an incorrect answer supported by substantially correct calculations and explanations might still receive partial credit.

- Answers should be provided in **English**.

# 1 UCB

Denote by $S_{j,t} = \sum_{k=1}^{t} X_{i_k,k} \cdot \mathbb{1}(i_k = a)$ and by $N_{j,t} = \sum_{k=1}^{t} \mathbb{1}(i_k = j)$ the cumulative reward and number of pulls of arm $j$ at time $t$. Denote by $\widehat{\mu}_{j,t} = \frac{S_{j,t}}{N_{j,t}}$ the estimated mean. Recall that, at each timestep $t$, UCB plays the arm $i_t$ such that

$$i_t \in \arg \max_j \widehat{\mu}_{j,t} + U(N_{j,t}, \delta)$$

Is $\widehat{\mu}_{j,t}$ an unbiased estimator (i.e., $\mathbb{E}_{UCB}[\widehat{\mu}_{j,t}] = \mu_j$)? Justify your answer.

# 2 Best Arm Identification

In best arm identification (BAI), the goal is to identify the best arm in as few samples as possible. We will focus on the fixed-confidence setting where the goal is to identify the best arm with high probability $1 - \delta$ in as few samples as possible. A player is given $k$ arms with expected reward $\mu_i$. At each timestep $t$, the player selects an arm to pull ($I_t$), and they observe some reward ($X_{I_t,t}$) for that sample. At any timestep, once the player is confident that they have identified the best arm, they may decide to stop.

**$\delta$-correctness and fixed-confidence objective.** Denote by $\tau_\delta$ the stopping time associated to the stopping rule, by $i^\star$ the best arm and by $\widehat{i}$ an estimate of the best arm. An algorithm is $\delta$-correct if it predicts the correct answer with probability at least $1 - \delta$. Formally, if $\mathbb{P}_{\mu_1,\ldots,\mu_k}(\widehat{i} \neq i^\star) \leq \delta$ and $\tau_\delta < \infty$ almost surely for any $\mu_1,\ldots,\mu_k$. Our goal is to find a $\delta$-correct algorithm that minimizes the sample complexity, that is, $\mathbb{E}[\tau_\delta]$ the expected number of sample needed to predict an answer.

Notation

- $I_t$: the arm chosen at round $t$.

- $X_{i,t} \in [0,1]$: reward observed for arm $i$ at round $t$.

- $\mu_i$: the expected reward of arm $i$.

- $\mu^\star = \max_i \mu_i$.

- $\Delta_i = \mu^\star - \mu_i$: suboptimality gap.

Consider the following algorithm

---
**Input:** $k$ arms, confidence $\delta$
$S = \{1, \ldots, k\}$
**for** $t = 1, \ldots$ **do**
$\quad$ Pull **all** arms in $S$
$\quad S = S \setminus \left\{ i \in S \ : \ \exists j \in S, \ \widehat{\mu}_{j,t} - U(t,\delta) \geq \widehat{\mu}_{i,t} + U(t,\delta) \right\}$
$\quad$ **if** $|S| = 1$ **then**
$\quad\quad$ STOP
$\quad\quad$ **return** $S$
$\quad$ **end**
**end**

---

The algorithm maintains an active set $S$ and an estimate of the empirical reward of each arm $\widehat{\mu}_{i,t} = \frac{1}{t} \sum_{j=1}^{t} X_{i,j}$.

- Compute the function $U(t, \delta)$ that satisfy the any-time confidence bound. For any arm $i \in [k]$

$$\mathbb{P}\left( \bigcup_{t=1}^{\infty} \{|\widehat{\mu}_{i,t} - \mu_i| > U(t,\delta)\} \right) \leq \delta$$

  Use Hoeffding's inequality.

- Let $\mathcal{E} = \bigcup_{i=1}^{k} \bigcup_{t=1}^{\infty} \{|\widehat{\mu}_{i,t} - \mu_i| > U(t,\delta')\}$. Using previous result shows that $\mathbb{P}(\mathcal{E}) \leq \delta$ for a particular choice of $\delta'$. This is called "bad event" since it means that the confidence intervals do not hold.

- Show that with probability at least $1 - \delta$, the optimal arm $i^\star = \arg\max_i\{\mu_i\}$ remains in the active set $S$. Use your definition of $\delta'$ and start from the condition for arm elimination. From this, use the definition of $\neg\mathcal{E}$.

- Under event $\neg\mathcal{E}$, show that an arm $i \neq i^\star$ will be removed from the active set when $\Delta_i \geq C_1 U(t, \delta')$ where $C_1 > 1$ is a constant. Compute the time required to have such condition for each non-optimal arm. Use the condition of arm elimination applied to arm $i^\star$.

- Compute a bound on the sample complexity (after how many rounds the algorithm stops) for identifying the optimal arm w.p. $1 - \delta$.

Note that also a variations of UCB are effective in pure exploration.

# 3  Bernoulli Bandits

In this exercise, you compare KL-UCB and UCB empirically with Bernoulli rewards $X_t \sim Bern(\mu_{I_t})$.

- Implement KL-UCB and UCB

**KL-UCB:**
$$I_t = \arg\max_i \max \left\{ \mu \in [0,1] : d(\widehat{\mu}_{i,t}, \mu) \leq \frac{\log(1 + t \log^2(t))}{N_{i,t}} \right\}$$

  where $d$ is the Kullback–Leibler divergence (see closed form for Bernoulli). A way of computing the inner max is through bisection (finding the zero of a function).

**UCB:**

$$I_t = \arg\max_i \widehat{\mu}_{i,t} + \sqrt{\frac{\log(1 + t\log^2(t))}{2N_{i,t}}}$$

that has been tuned for $1/2$-subgaussian problems.

- Let $n = 10000$ and $k = 2$. Plot the <u>expected</u> regret of each algorithm as a function of $\Delta$ when $\mu_1 = 1/2$ and $\mu_2 = 1/2 + \Delta$.

- Repeat the above experiment with $\mu_1 = 1/10$ and $\mu_1 = 9/10$.

- Discuss your results.

# 4 Regret Minimization in RL

Consider a finite-horizon MDP $M^\star = (S, A, p_h, r_h)$ with stage-dependent transitions and rewards. Assume rewards are bounded in $[0, 1]$. We want to prove a regret upper-bound for UCBVI. We will aim for the suboptimal regret bound ($T = KH$)

$$R(T) = \sum_{k=1}^{K} V_1^\star(s_{1,k}) - V_1^{\pi_k}(s_{1,k}) = \widetilde{O}(H^2 S\sqrt{AK})$$

Define the set of plausible MDPs as

$$\mathcal{M}_k = \{M = (S, A, p_{h,k}, r_{h,k}) \ : \ r_{h,k}(s, a) \in \beta_{h,k}^r(s, a), p_{h,k}(\cdot|s, a) \in \beta_{h,k}^p(s, a)\}$$

Confidence intervals can be anytime or not.

- Define the event $\mathcal{E} = \{\forall k, M^\star \in \mathcal{M}_k\}$. Prove that $\mathbb{P}(\neg\mathcal{E}) \leq \delta/2$. First step, construct a confidence interval for rewards and transitions for each $(s, a)$ using Hoeffding and Weissmain inequality (see appendix), respectively. So, we want that

$$\mathbb{P}\Big(\forall k, h, s, a : |r_{hk}(s, a) - r_h(s, a)| \leq \beta_{hk}^r(s, a) \wedge \|\widehat{p}_{hk}(\cdot|s, a) - p_h(\cdot|s, a)\|_1 \leq \beta_{hk}^p(s, a)\Big) \geq 1 - \delta/2$$

- Define the bonus function and consider the Q-function computed at episode $k$

$$Q_{h,k}(s, a) = \widehat{r}_{h,k}(s, a) + b_{h,k}(s, a) + \sum_{s'} \widehat{p}_{h,k}(s'|s, a)V_{h+1,k}(s')$$

with $V_{h,k}(s) = \min\{H, \max_a Q_{h,k}(s, a)\}$. Recall that $V_{H+1,k}(s) = V_{H+1}^\star(s) = 0$. Prove that under event $\mathcal{E}$, $Q_k$ is optimistic, i.e.,

$$Q_{h,k}(s, a) \geq Q_h^\star(s, a), \forall s, a$$

where $Q^\star$ is the optimal Q-function of the unknown MDP $M^\star$. Note that $\widehat{r}_{H,k}(s, a) + b_{h,k}(s, a) \geq r_{h,k}(s, a)$ and thus $Q_{H,k}(s, a) \geq Q_H^\star(s, a)$ (for a properly defined bonus). Then use induction to prove that this holds for all the stages $h$.

- In class we have seen that

$$\delta_{hk}(s_{1,k}) \leq \sum_{h=1}^{H} Q_{hk}(s_{hk}, a_{hk}) - r(s_{hk}, a_{hk}) - \mathbb{E}_{Y \sim p(\cdot|s_{hk}, a_{hk})}[V_{h+1,k}(Y)]) + m_{hk} \quad (1)$$

where $\delta_{hk}(s) = V_{hk}(s) - V_h^{\pi_k}(s)$ and $m_{hk} = \mathbb{E}_{Y \sim p(\cdot|s_{hk}, a_{hk})}[\delta_{h+1,k}(Y)] - \delta_{h+1,k}(s_{h+1,k})$. We now want to prove this result. Denote by $a_{hk}$ the action played by the algorithm (you will have to use the greedy property).

1. Show that $V_h^{\pi_k}(s_{hk}) = r(s_{hk}, a_{hk}) + \mathbb{E}_p[V_{h+1,k}(s')] - \delta_{h+1,k}(s_{h+1,k}) - m_{h,k}$

Initialize $Q_{h1}(s,a) = 0$ for all $(s,a) \in S \times A$ and $h = 1, \dots, H$

**for** $k = 1, \dots, K$ **do**
    Observe initial state $s_{1k}$ *(arbitrary)*
    Estimate empirical MDP $\widehat{M}_k = (S, A, \widehat{p}_{hk}, \widehat{r}_{hk}, H)$ from $\mathcal{D}_k$

$$\widehat{p}_{hk}(s'|s,a) = \frac{\sum_{i=1}^{k-1} \mathbb{1}\{(s_{hi}, a_{hi}, s_{h+1,i}) = (s, a, s')\}}{N_{hk}(s,a)}, \quad \widehat{r}_{hk}(s,a) = \frac{\sum_{i=1}^{k-1} r_{hi} \cdot \mathbb{1}\{(s_{hi}, a_{hi}) = (s, a)\}}{N_{hk}(s,a)}$$

    Planning (by backward induction) for $\pi_{hk}$ using $\widehat{M}_k$
    **for** $h = H, \dots, 1$ **do**
       $Q_{h,k}(s,a) = \widehat{r}_{h,k}(s,a) + b_{h,k}(s,a) + \sum_{s'} \widehat{p}_{h,k}(s'|s,a) V_{h+1,k}(s')$
       $V_{h,k}(s) = \min\{H, \max_a Q_{h,k}(s,a)\}$
    **end**
    Define $\pi_{h,k}(s) = \arg\max_a Q_{h,k}(s,a), \forall s, h$
    **for** $h = 1, \dots, H$ **do**
       Execute $a_{hk} = \pi_{hk}(s_{hk})$
       Observe $r_{hk}$ and $s_{h+1,k}$
       $N_{h,k+1}(s_{hk}, a_{hk}) = N_{h,k}(s_{hk}, a_{hk}) + 1$
    **end**
**end**

**Algorithm 1:** UCBVI

2. Show that $V_{h,k}(s_{hk}) \leq Q_{h,k}(s_{hk}, a_{hk})$.

3. Putting everything together prove Eq. 1.

- Since $(m_{hk})_{hk}$ is an MDS, using Azuma-Hoeffding we show that with probability at least $1 - \delta/2$

$$\sum_{k,h} m_{hk} \leq 2H \sqrt{KH \log(2/\delta)}$$

Show that the regret is upper bounded with probability $1 - \delta$ by

$$R(T) \leq \sum_{kh} b_{hk}(s_{hk}, a_{hk}) + 2H \sqrt{KH \log(2/\delta)}$$

- Finally, we have that

$$\sum_{h,k} \frac{1}{\sqrt{N_{hk}(s_{hk}, a_{hk})}} = \sum_{h=1}^{H} \sum_{s,a} \sum_{i=1}^{N_{h,K}(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=1}^{H} \sum_{s,a} \sqrt{N_{hK}(s,a)}$$

Complete this by showing an upper-bound of $H\sqrt{SAK}$, which leads to $R(T) \lesssim H^2 S \sqrt{AK}$

# A   Weissmain inequality

Denote by $\widehat{p}(\cdot|s,a)$ the estimated transition probability build using $n$ samples drawn from $p(\cdot|s,a)$. Then we have that

$$\mathbb{P}(\|\widehat{p}_h(\cdot|s,a) - p_h(\cdot|s,a)\|_1 \geq \epsilon) \leq (2^S - 2) \exp\left(-\frac{n\epsilon^2}{2}\right)$$

# References