# Project: Performance of Multi-armed Bandit Algorithms

## Due Date: 11:59am, Dec. 31, 2019

- **Important Issue**

  You are required to use Python or R for the programming part. You also need to submit a report including your simulation results, analysis, discussions, tables and figures. You can obtain 10% bonus points if you submit the report and Python code together by using the Jupyter Notebook.

- **Basic Setting**

  We consider a time-slotted bandit system ($t = 0, 1, 2, \dots$) with three arms. We denote the arm set as $\{1, 2, 3\}$. Pulling each arm $j$ ($j \in \{1, 2, 3\}$) will obtain a reward $r_j$, which satisfies a Bernoulli distribution with mean $\theta_j$ ($\mathrm{Bern}(\theta_j)$), $i.e.$,

  $$r_j = \begin{cases} 1, & w.p. \ \theta_j, \\ 0, & w.p. \ 1 - \theta_j, \end{cases}$$

  where $\theta_j$ are parameters within $(0, 1)$ for $j \in \{1, 2, 3\}$.

  Now we run this bandit system for $N$ ($N \gg 3$) time slots. At each time slot $t$, we choose one and only one arm from these three arms, which we denote as $I(t) \in \{1, 2, 3\}$. Then we pull the arm $I(t)$ and obtain a reward $r_{I(t)}$. Our objective is to find an optimal policy to choose an arm $I(t)$ at each time slot $t$ such that the expectation of the aggregated reward is maximized, $i.e.$,

  $$\max_{I(t), t=1, \dots, N} \ \mathbb{E}\left[\sum_{t=1}^{N} r_{I(t)}\right].$$

  If we know the values of $\theta_j, j \in \{1, 2, 3\}$, this problem is trivial. Since $r_{I(t)} \sim \mathrm{Bern}(\theta_{I(t)})$,

  $$\mathbb{E}\left[\sum_{t=1}^{N} r_{I(t)}\right] = \sum_{t=1}^{N} \mathbb{E}[r_{I(t)}] = \sum_{t=1}^{N} \theta_{I(t)}.$$

  Let $I(t) = I^* = \arg\max_{j \in \{1,2,3\}} \theta_j$ for $t = 1, 2, \dots, N$, then

  $$\max_{I(t), t=1, \dots, N} \ \mathbb{E}\left[\sum_{t=1}^{N} r_{I(t)}\right] = N \cdot \theta_{I^*}.$$

  However, in reality, we do not know the values of $\theta_j, j \in \{1, 2, 3\}$. We need to estimate the values $\theta_j$ via empirical samples, and then make the decisions at each time slot.

  Next we introduce three classical bandit algorithms: $\epsilon-$greedy, UCB and Thompson sampling.

- **Bandit Algorithms**

  1. $\epsilon-$greedy Algorithm ($0 \le \epsilon \le 1$)

---

**Algorithm 1** $\epsilon-$greedy Algorithm

---

**Initialize** $\hat{\theta}(j) = 0, \text{count}(j) = 0, j \in \{1, 2, 3\}$

1: **for** $t = 1, 2 \ldots, N$ **do**

2:

$$I(t) \leftarrow \begin{cases} \underset{j \in \{1,2,3\}}{\arg\max} \ \hat{\theta}(j) & w.p. \ 1 - \epsilon \\ \text{randomly chosen from}\{1,2,3\} & w.p. \ \epsilon \end{cases}$$

3: $\quad$ $\text{count}(I(t)) \leftarrow \text{count}(I(t)) + 1$

4: $\quad$ $\hat{\theta}(I(t)) \leftarrow \hat{\theta}(I(t)) + \frac{1}{\text{count}(I(t))} \left[ r_{I(t)} - \hat{\theta}(I(t)) \right]$

5: **end for**

---

2. UCB (Upper Confidence Bound) Algorithm

---

**Algorithm 2** UCB Algorithm

---

1: **for** $t = 1, 2, 3$ **do**

2: $\quad$ $I(t) \leftarrow t$

3: $\quad$ $\text{count}(I(t)) \leftarrow 1$

4: $\quad$ $\hat{\theta}(I(t)) \leftarrow r_{I(t)}$

5: **end for**

6: **for** $t = 4, \ldots, N$ **do**

7:

$$I(t) \leftarrow \underset{j \in \{1,2,3\}}{\arg\max} \left( \hat{\theta}(j) + c \cdot \sqrt{\frac{2 \log t}{\text{count}(j)}} \right)$$

8: $\quad$ $\text{count}(I(t)) \leftarrow \text{count}(I(t)) + 1$

9: $\quad$ $\hat{\theta}(I(t)) \leftarrow \hat{\theta}(I(t)) + \frac{1}{\text{count}(I(t))} \left[ r_{I(t)} - \hat{\theta}(I(t)) \right]$

10: **end for**

---

**Note**: $c$ is a positive constant with a default value of 1.

3. Thompson sampling (TS) Algorithm

Recall that $\theta_j, j \in \{1, 2, 3\}$, are unknown parameters over $(0, 1)$. From the Bayesian perspective, we assume their priors are Beta distributions with given parameters $(\alpha_j, \beta_j)$.

---
**Algorithm 3** Thompson sampling Algorithm
---
**Initialize**  Beta parameter $(\alpha_j, \beta_j), j \in \{1, 2, 3\}$

1: **for** $t = 1, 2 \ldots, N$ **do**
2:    *# Sample model*
3:    **for** $j \in \{1, 2, 3\}$ **do**
4:       Sample $\hat{\theta}(j) \sim \text{Beta}(\alpha_j, \beta_j)$
5:    **end for**
6:    *# Select and pull the arm*

$$I(t) \leftarrow \arg\max_{j \in \{1,2,3\}} \hat{\theta}(j)$$

7:    *# Update the distribution*

$$\alpha_{I(t)} \leftarrow \alpha_{I(t)} + r_{I(t)}$$
$$\beta_{I(t)} \leftarrow \beta_{I(t)} + 1 - r_{I(t)}$$

8: **end for**
---

- **Simulation**

  1. Suppose we obtain the Bernoulli distribution parameters from an oracle, which are shown in the following table below. Choose $N = 10000$ and compute the theoretically maximized expectation of aggregate rewards over $N$ time slots. We call it the oracle value. Note that these parameters $\theta_j, j = 1, 2, 3$ and oracle values are unknown to the three bandit algorithms.

     | Arm $j$ | 1 | 2 | 3 |
     |---------|-----|-----|-----|
     | $\theta_j$ | 0.4 | 0.6 | 0.8 |

  2. Implement three bandit algorithms with following settings:
     - $\epsilon$-greedy with $\epsilon = 0.1, 0.5, 0.9$.
     - UCB with $c = 1, 5, 10$.
     - Thompson Sampling with
       $\{(\alpha_1, \beta_1) = (1, 1), (\alpha_2, \beta_2) = (1, 1), (\alpha_3, \beta_3) = (1, 1)\}$ and
       $\{(\alpha_1, \beta_1) = (2, 4), (\alpha_2, \beta_2) = (3, 6), (\alpha_3, \beta_3) = (1, 2)\}$

  3. For each algorithm, repeat the experiment 100 times and use the average of results of 100 experiments.

  4. Store the reward in each time slot $t = 1, 2, 3, ..., 10000$. For Thompson Sampling, also store the (posterior) probability of playing each arm.

  5. Plot the reward and the average reward against time slot ($t = 1, 2, 3..., N$). For Thompson Sampling, plot the arm playing probabilities as well.

  6. Compute the gaps between the algorithm outputs and the oracle value. Compare the numerical results of $\epsilon$-greedy, UCB, and Thompson Sampling, Which one is the best? Then discuss the impacts of $\epsilon$, $C$, and $\alpha_j$, $\beta_j$, respectively.

  7. Give your understanding of the exploration-exploitation trade-off in bandit algorithms.

  8. We implicitly assume the reward distribution of three arms are independent. How about the dependent case? Can you design an algorithm to exploit such information to obtain a better result?