

9 Multi-Armed Bandit

Exercise 9.1

Tell if the following statements about MAB are true or false and provide motivations for your choice.

1. The MDP corresponding to a MAB setting has no state.
2. The MDP corresponding to a MAB setting has no transition probabilities.
3. An algorithm which always uses the greedy choice (maximize the empirical expected reward) might get stuck to suboptimal solutions.
4. We are able to solve the exploration/exploitation problem by considering either upper or lower bounds on the expected rewards.
5. In a frequentist framework, if we base our sequential policy on a MAB by considering only $\hat{R}(a_i)$ we are not using all the information we have about the estimates.
6. Uncertainty over quantities is usually handled by explorative choices in the MAB setting.

Exercise 9.2

Tell if the following problems can be modeled as MAB and explain why.

1. Maze escape;
2. Pricing goods with stock;
3. Pricing goods without stock;
4. Optimal bandwidth allocation (send the largest amount of information without congesting the band);
5. Web ads placement;
6. Weather prediction (with multiple experts).

If they fit the MAB setting, is the environment adversarial or stochastic?

Exercise 9.3

Provide an example for which the pure exploitation strategy is failing to converge to the optimum in a MAB setting with Bernoulli rewards. Recall that the pure exploitation algorithm chooses the arm by selecting the one with the largest:

$$\hat{R}_t(a_i) = \frac{1}{N_t(a_i)} \sum_{j=1}^t r_{i,j} \mathbb{I}\{a_i = a_{i_j}\},$$

where \mathbb{I} is the indicator function.

How often does this occurs? Can we compute an lower bound over the regret of this strategy?

Exercise 9.4

The ε -greedy algorithm selects the best action except for small percentages of times ε , where all the actions are considered. Consider a MAB stochastic setting.

1. Is this algorithm converging to the optimal strategy (in some sense)?
2. If not, propose a scheme which has the chance of converging to the optimal solution.
3. Are we in a MAB perspective if we are using this algorithm?

Exercise 9.5

Write the formula for the minimum regret we might have on average over $T = \lceil e^{10} \rceil$ time steps in the case we have a stochastic MAB with 3 arms and expected rewards:

$$R(a_1) = 0.2 \tag{9.1}$$

$$R(a_2) = 0.4 \tag{9.2}$$

$$R(a_3) = 0.7 \tag{9.3}$$

and each distribution $\mathcal{R}(a_i)$ is Bernoulli.

Note that the KL divergence for Bernoulli variables with means p and q is:

$$KL(p, q) = p \log \left(\frac{p}{q} \right) + (1 - p) \log \left(\frac{(1 - p)}{(1 - q)} \right).$$

Hints: $\log(\frac{0.2}{0.7}) = -1.25$, $\log(\frac{0.8}{0.3}) = 0.98$, $\log(\frac{0.4}{0.7}) = -0.56$ and $\log(\frac{0.6}{0.3}) = 0.69$.

Is it possible that your algorithm achieves lower regret? If so, provide an example.

Exercise 9.6

Provide examples of either Bayesian or frequentist MAB algorithm showing the following properties:

1. It incorporates expert knowledge about the problem in the arms;
2. It provides tight theoretical lower bound on the expected regret in the stochastic setting;
3. It provides tight theoretical upper bound on the expected regret in the stochastic setting;
4. At each turn, it modifies only the statistics of the chosen arm.

Motivate your answers.

Exercise 9.7

Consider the following bounds (supposed to hold with probability at least δ) for a MAB stochastic setting:

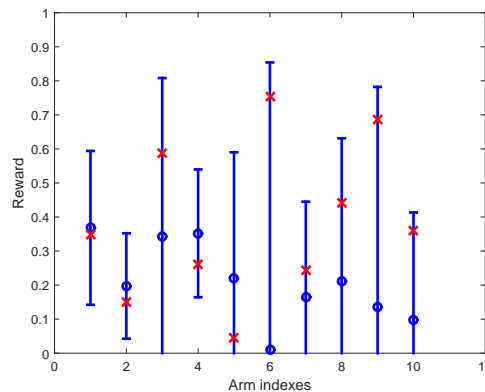


Figure 9.1: Bounds in a MAB with $N = 10$.

where the blue bars are the Hoeffding bounds for the expected rewards, the blue circles are the estimated expected rewards, and the red crosses are the real expected rewards.

1. Which arm would a UCB1 algorithm choose for the next round?
2. Do you think that Figure 9.1 might be the results obtained by running UCB1 for several rounds?

3. Which arm will UCB1 converge to if $T \rightarrow \infty$?
4. Which arm is the one which we pulled the most so far?

Motivate your answers.

Exercise 9.8

Consider the Thompson Sampling algorithm. Assume to have the following posterior distributions $Beta_i(\alpha_t, \beta_t)$ for arms $\mathcal{A} = \{a_1, \dots, a_5\}$ rewards, which are distributed as Bernoulli r.v.:

$a_1 :$	$\alpha_t = 1$	$\beta_t = 5$
$a_2 :$	$\alpha_t = 6$	$\beta_t = 4$
$a_3 :$	$\alpha_t = 11$	$\beta_t = 23$
$a_4 :$	$\alpha_t = 12$	$\beta_t = 33$
$a_5 :$	$\alpha_t = 28$	$\beta_t = 21$

From these distribution you extract the following samples for the current round:

$$\begin{aligned}\hat{r}(a_1) &= 0.63 \\ \hat{r}(a_2) &= 0.35 \\ \hat{r}(a_3) &= 0.16 \\ \hat{r}(a_4) &= 0.22 \\ \hat{r}(a_5) &= 0.7\end{aligned}$$

1. Which arm would the TS algorithm play for the next round?
2. What changes if the real distributions of the arm rewards are Gaussian.
3. Assume we started the TS algorithm with uniform $Beta(1, 1)$ priors. What would UCB1 have chosen in the case of Bernoulli rewards for the next round?

Solutions

Answer of exercise 9.1

1. FALSE It has a single state.
2. FALSE The transition probability matrix is $P(s, a_i | s) = 1$.
3. TRUE There are cases in which if we only consider the expected values of the rewards as decision tool, we might discard good options even if we only gathered small evidence from them.
4. FALSE Pessimistic choices would not solve the exploration/exploitation dilemma since, this way, we do not provide incentives to explore the options which are considered suboptimal so far.
5. TRUE We are also allowed to consider concentration inequalities which provides us information about the uncertainty we have about the reward estimate.
6. TRUE Sometimes, we want to reduce uncertainty by picking suboptimal choices which allow us to gather more information that will be beneficial in the future.

Answer of exercise 9.2

1. NO There are plenty of states in a maze escape problem (the reward cannot be modeled as a single distribution).
2. NO We are limited in the number of successes, thus we might stop the problem at some point (it can be a generalized MAB, in the specific a budget MAB).
3. YES Rewards are stochastic if we consider a distribution of users or adversarial if we assume to have a single malicious buyer.
4. YES If we consider the traffic in the bandwidth as a stochastic event.
5. YES Assuming, for instance, a pay-per-click scheme and a click event which is a stochastic event.
6. NO Each time we predict we have rewards coming from all the arms (expert problem). If we use only the feedback coming from the pulled arm, we can use the MAB model, even if we would discard some useful information by doing this.

Answer of exercise 9.3

The most simple example is the one in which the first time we pull the optimal arm

we have a zero reward. This occurs with probability $1 - \mu^*$ and thus, being positive probability that the arm is pulled only once, the regret cannot be less than linear in the time horizon, i.e., $O(T)$.

Answer of exercise 9.4

1. No, it can not reach the optimal policy since it will always select the suboptimal arms with a fixed probability. This leads to a regret of the order $\varepsilon T \gg O(\log T)$.
2. If we use a strategy s.t. $\lim_{t \rightarrow \infty} \varepsilon(t) = 0$ we might converge.
3. Even if this strategy converges, you do not have any assurance about the regret we are suffering in the process of learning. There exists a theoretical analysis of a version of the ε -greedy algorithm in:
Auer, Peter, Nicolo Cesa-Bianchi, and Paul Fischer. "Finite-time analysis of the multi-armed bandit problem." Machine learning 47.2-3 (2002): 235–256
 The solution proposed can be applied only by knowing the gaps Δ_i between arms, otherwise the bound might not hold.

Answer of exercise 9.5

If we assume that T is sufficiently large, we have:

$$R_T \geq \log T \sum_{a_i \neq a^*} \frac{\Delta_i}{KL(\mathcal{R}(a_i), \mathcal{R}(a^*))},$$

thus in our case:

$$\begin{aligned} R_T &\geq \log e^{10} \left(\frac{0.7 - 0.2}{KL(0.2, 0.7)} + \frac{0.7 - 0.4}{KL(0.4, 0.7)} \right) \\ &= 10 \left(\frac{0.5}{0.2(-1.25) + 0.8(0.98)} + \frac{0.3}{0.4(-0.56) + 0.6(0.69)} \right) \\ &= 25.1528 \end{aligned}$$

We might violate our lower bound in some cases, e.g.:

- if we consider a subset of MAB settings;
- if we consider a single realization of the rewards.

Answer of exercise 9.6

1. TS As usual we can use prior distributions to incorporate information about the problem.

2. NONE The lower bound is defined over the problem and not by basing on the algorithm we used.
3. UCB1 and TS Provide order optimal theoretical upper bound on the expected regret in the stochastic setting. TS matches also the constant (KL-UCB has the same assurance).
4. TS It only modifies the Beta distribution corresponding to the pulled arm.

Answer of exercise 9.7

1. a_6 since it is the one with the highest upper bound.
2. No, since in principle UCB1 keeps all the upper bounds at similar level.
3. a_6 , since it is the one with highest expected value.
4. a_2 , since it is the one with smallest bounds (which are inversely proportional to the number of pulls).

Answer of exercise 9.8

1. Thompson sampling will choose the arm a_i with the largest sampled $\hat{r}(a_i)$, which in this case is a_5 .
2. In the case we are considering Gaussian rewards it is not meaningful to use Beta distribution as prior. For instance, if the reward are Gaussian they might provide also negative values, while the support of the Beta is $[0, 1]$. Thus, it makes no sense to instantiate a problem with Beta prior when you have Gaussian rewards.
3. Since we started with uniform prior and at each round we collected a success or a failure we are currently at round:

$$t = 4 + 8 + 32 + 43 + 47 = 134,$$

while the UCB1 upper bound $U_t(a_i)$ is of the form:

$$U_t(a_i) = \frac{\alpha_t - 1}{\alpha_t + \beta_t - 2} + \sqrt{\frac{2 \log t}{\alpha_t + \beta_t - 2}},$$

thus:

$$\begin{aligned}U_t(a_1) &= \frac{0}{4} + \sqrt{\frac{2 \log 134}{4}} \\U_t(a_2) &= \frac{5}{8} + \sqrt{\frac{2 \log 134}{8}} \\U_t(a_3) &= \frac{10}{32} + \sqrt{\frac{2 \log 134}{32}} \\U_t(a_4) &= \frac{11}{43} + \sqrt{\frac{2 \log 134}{43}} \\U_t(a_5) &= \frac{27}{47} + \sqrt{\frac{2 \log 134}{47}}\end{aligned}$$

In this case the UCB1 algorithm chooses the arm a_i s.t. $i = \arg \max_i U_t(a_i)$, in this case a_2 .