

Introduction to Contextual Bandits

Instructors: Alekh Agarwal, Sham Kakade

Lecture 6

6.1 Formal Setting and Examples

The contextual bandit (CB) setting is best described as a repeated interaction which proceeds over T rounds.

At each round $t = 1, 2, \dots, T$:

1. Environment reveals a context $x_t \in \mathcal{X}$.
2. Learner chooses an action $a_t \in \mathcal{A}$.
3. Environment reveals a reward $r_t \in [0, 1]$.

Figure 6.1: The repeated interaction between the environment and a learning agent in the contextual bandit problem

The goal of the learner is to choose actions which maximize its cumulative reward, that is $\sum_{t=1}^T r_t$. Throughout, we will assume that there is an (unknown) joint distribution over contexts as well as the rewards of all the actions, denoted as D_t at time t . With some abuse of notation, we will also use D_t to denote the marginal distribution over contexts and conditional distribution over the rewards given a context. With this, the above interaction involves environment choosing $x_t \sim D_t$ and $r_t \sim D_t(\cdot | x_t, a_t)$. The learner is equipped with a set of policies $\Pi \subseteq \{\mathcal{X} \rightarrow \mathcal{A}\}$, and aims to find a policy $\pi \in \Pi$ which attains a large reward. The (expected) average reward of a policy π , often also called the value of π is defined as

$$V(\pi) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{x \sim D_t} [\mathbb{E}[r_t | x, \pi(x)]] . \quad (6.1)$$

We will often use the shorthand $\mathbb{E}_{D_t}[r(\pi(x))]$ to refer to the above expectation concisely.

We will focus on a setting where the set Π is either extremely large, or infinite. Similarly, we will typically consider the set \mathcal{X} to be infinite, so that learner has to *generalize* across contexts. The set of actions \mathcal{A} , on the other hand will be expected to have a reasonable cardinality, with $|\mathcal{A}|$ denoted by K . Indeed, we will often identify the actions set with $[K]$, the set of integers $\{1, 2, \dots, K\}$.

With the formal notation in place, we now flesh out a couple of examples.

Example 1 (Online news recommendation). Let us consider a news website, which has a pool of candidate stories at any given point of time. When a user visits the website, it wants to decide which stories to display to the user. To simplify the discussion, let us assume for now that only one story out of the pool will be actually presented to the user. In this case, the context x can consist of any information about the user, such as their location, browser, any reading habits as gleaned from prior visits to the website, etc. The context might further contain features describing the content of the various news stories in the pool. The action set \mathcal{A} consists of the pool of news stories, with the tacit assumption that this set is not too large, say at most 30. Rewards should measure user engagement with the chosen news story, which can be measured in a number of ways. We could use clicks, meaning a reward of 1 if the chosen story is clicked on by the user. More fine-grained rewards could further encode the time spent by the user reading the story, whether they explicitly like or share the story on social media etc. Policies in this setting might extract some features from

the available information, and predict a real number for every action describing its quality. Concretely, if the features are described as a vector in $x \in \mathbb{R}^d$, then a policy might maintain a weight vector $w_i \in \mathbb{R}^d$ for each action i , and the action picked by the policy can be $\arg \max_{i=1}^K w_i^T x$ with ties broken in favor of smaller i . More generally, we can use any other multiclass classifier, such as based on tree ensembles, neural nets, etc. and this forms the most natural policy sets in practice.

Notice that we picked reward measures which can be clearly associated with a particular action, such as a click or time spent reading the selected story above. Often, longer term measures such as the number of visits by a user in a month might be more accurate measures of user engagement with the website. However, maximizing such metrics falls outside the scope of the contextual bandit problem, and results in harder reinforcement learning problems which we will visit in the later part of the course. Often the best practical approach is to find a well-defined short-term metric such as time spent on a story, which might be well correlated with desired longer term measures.

Example 2 (Mobile health advice). As our second example, we consider a mobile health application which is designed to prompt people suffering from a chronic condition such as diabetes or obesity to take a short exercise break. Let us assume, for simplicity, that the app displays these messages at 3 fixed times in the day so that the actions do not involve choosing the time of messaging. Suppose we have a small set of candidate messages, but not all phrasings and not all exercise options succeed equally for all users. Furthermore, the ideal choice might be different during the day when a person is in their office, and in the evening when they are at home. In this case, the context consists of broad health information about a person such as gender, age, height, weight etc., as well as some coarse summary of their previous interactions with the app (such as what type of workouts they tend to prefer). Actions are the different messages we can display and policies can be designed along similar lines as the previous example. A reward might be a binary indicator of whether the suggestion was taken or not, either indicated by the user explicitly, or measured using a wearable device.

Note that in both the above examples, part of the context is assumed to be based on the previous interactions between the agent and its environment, through past reading habits or the exercise suggestions which were followed. This is typical of many contextual bandit problems, where even in a non-adversarial setting, the future contexts are affected by the chosen actions. However, this effect is often relatively small over short and medium time-spans and hence we choose to ignore the different distributions that some other policy might have generated when we evaluate it.

We conclude this section with one final example.

Example 3 (Faulty VM handling). In this example, consider the problem of dealing with faulty virtual machines (VMs) in a large cloud service. The context consists of various telemetry readings about the workload on the VM when it stopped responding, hardware and software profile of the physical machine and the VM as well as broad indicators of the general health of the data center or the rack where the VM is located. Actions can be discrete waiting times, say 1 second apart between a smallest and largest reasonable duration, before which the VM should be either rebooted or migrated. If the machine recovers during the chosen waiting period, then the reward can be that chosen period. If not, the reward can be the chosen waiting time plus the time cost of a reboot or a migration.

Finally, we describe the objective of learning in contextual bandit problems more concretely. A typical learning objective is to minimize regret, which is defined as

$$\text{Regret}_T = \max_{\pi \in \Pi} \sum_{t=1}^T \mathbb{E}_{D_t} [r(\pi(x))] - \sum_{t=1}^T \mathbb{E}_{D_t} [r(a_t)]. \quad (6.2)$$

The definition looks rather complicated, and it is worthwhile unpacking it a little bit for two versions of the problem, i.i.d. and adversarial as shown in Figure 6.1. In the i.i.d. setting, the distributions D_t are identical for all rounds t . In the adversarial setting, the rewards and contexts at round t can depend on all the previous actions of the the player. However, it turns out that the worst-case adversary does not need to choose contexts and rewards from a distribution.

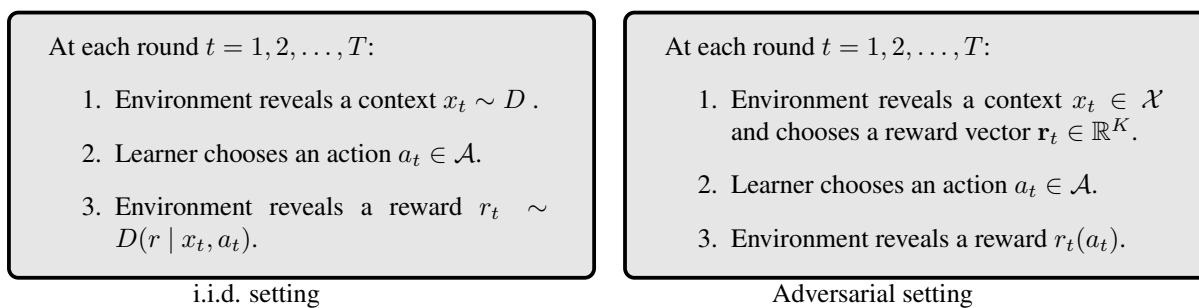


Figure 6.2: The repeated interaction between the environment and a learning agent for stochastic and adversarial CB

That is, the worst-case strategy for the adversary uses degenerate distributions D_t , since the rewards are only revealed to the player after the action a_t has been taken. Consequently, we can also specify the notions of regret for these settings as

$$\text{Regret}_T = \max_{\pi \in \Pi} \sum_{t=1}^T \mathbb{E}_D[r(\pi(x_t))] - \sum_{t=1}^T \mathbb{E}_D[r(a_t)]. \quad \text{i.i.d. setting} \quad (6.3)$$

$$\text{Regret}_T = \max_{\pi \in \Pi} \sum_{t=1}^T r_t(\pi(x_t)) - \sum_{t=1}^T r_t. \quad \text{adversarial setting} \quad (6.4)$$

6.2 Connections to other learning settings

In this section, how we will review how contextual bandits relate to multi-armed bandits, supervised learning and reinforcement learning.

Connections to multi-armed bandits: The similarities and differences between multi-armed and contextual bandits are most apparent from the notion of regret (6.2). While we also consider regret in multi-armed bandits (MAB), the algorithm's cumulative reward is compared to that of the *best fixed action in hindsight*. This corresponds to an extremely simple policy class of ignoring the context completely. Clearly this makes MAB a special case of CB. However, this special case can be overly restricted in most practical applications. Considering our earlier examples, applying MAB to news recommendation, mobile health and faulty VM handling corresponds to making the same choice, irrespective of who the user is or the characteristics of the VM in question. Contextual bandits with an expressive policy class can allow us to obtain a much higher reward in most of these settings.

Connections to reinforcement learning: At the other extreme, reinforcement learning (RL) is a strict generalization of CB. Loosely stated, while the agent's actions influence the future state of the world in RL, this is not the case in CB. More formally, we again examine the notion of regret in adversarial CB problems (6.4). In this case, the contexts and rewards chosen by the adversary at time t depend on the actions of the algorithm at all the past rounds. That is, the algorithm's choices affect the future data it experiences. Crucially though, the algorithm and any other policy π in the second term of regret are both evaluated under the same contexts and rewards. We do not account for the fact that if the algorithm had indeed followed the policy π , then the adversary would have responded with different contexts and rewards. In RL, we would indeed consider this more complex benchmark (called policy regret), where each policy is evaluated on the data distribution that would be generated in response to its actions. Thus, CB makes a simplifying assumption compared to RL.

It turns out this assumption is already partly violated even in the examples we discussed earlier. In news or health advice recommendation tasks, it is common to use features which rely on the previous history of a user's interaction, such as the number of times they took a particular suggestion. Of course, this depends on the policy making these recommendations and formally violate the CB assumption. As we will see in the sequel, this simplification is critical for the development of practical algorithms, and still ends up being a reasonable approximation for many problems of interest. We also notice that the simplification is automatically correct in the i.i.d. version of the problem.

Connections to supervised learning: In practice, many CB problems are often cast as supervised learning problems instead. For instance, suppose we are interested in news recommendation. Then for every user, we record the chosen article and whether they read it or not. We now cast this as a conditional probability prediction problem, where given a (user, article) pair, our goal is to predict the probability of the user clicking on the article. Based on the labels collected in our system logs, we can easily run a logistic or least squares regression and obtain a model for predicting these probabilities. Seems simple enough, why use CB?

The problem with the above approach is that our logistic regression model is only guaranteed to be accurate on the distribution of data which it is trained on. Suppose we have an existing system A that is used to recommend articles to users. Then we know that for the distribution of (user, article) pairs generated by system A, we have a good estimate of click probabilities. Once we have such a model though, we want to use it for selecting the articles that are displayed, presumably based on the largest predicted click probability. Furthermore, this results in the choice of very different articles than the ones picked by system A. In an extreme version, imagine that system A just picks the same article for every user. Then we have no statistical evidence for or against the quality of any other article, and our logistic regression model is completely blind to their quality! Consequently, the training and test errors we will compute based on the data collected under system A will have little connection with the number of clicks we will see online when we deploy our logistic regression model.

Once again, we can contrast with the CB model by examining the notion of regret. In regret, we compare the cumulative reward of our choices, with that of choices made according to any other policy. Consequently, a small regret directly informs us that no other policy in our class can allow us to obtain a substantially higher reward, unlike the test error of logistic regression model above. In the next lecture, we will also see a suited of techniques that set up a more formal training/test style paradigm for contextual bandit problems as an alternative to regret.

6.3 Warm Up: Stochastic Multi-Armed Bandits

We begin our discussion with the most classical example from this class of problems, where there is no context, and the rewards of each action are drawn i.i.d. from a distribution at each round. Formally, arm a has an associated distribution D_a and $r_t(a) \sim D_a$ for all $t = 1, 2, \dots, T$. Let $\mu_a = \mathbb{E}_{D_a}[r]$ and without loss of generality, we assume that r_1, r_2, \dots, r_K are described in the order of decreasing μ_a . We also define $\Delta_a = \mu_1 - \mu_a$. We now describe an extremely intuitive strategy called the upper confidence bound (UCB) strategy for stochastic multi-armed bandits. In this strategy, let $n_t(a)$ refer to the number of times arm a has been pulled prior to round t , and let $\hat{\mu}_t(a)$ be the sample average of the rewards obtained by action a on the rounds where it was played. Then at round t , UCB takes as input a probability parameter δ and chooses the action

$$i_t = \operatorname{argmax}_{a=1,2,\dots,K} \hat{\mu}_t(a) + \sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_t(a)}}. \quad (6.5)$$

We will show the following result regarding the regret of the UCB strategy.

Theorem 4. *With probability at least $1 - \delta$, the regret of the UCB strategy is at most:*

$$T\mu(1) - \sum_{t=1}^T \mu(i_t) \leq K - 1 + \sum_{a=2}^K \frac{2 \ln \frac{2KT}{\delta}}{\Delta_a}.$$

Proof: Let us first consider the event \mathcal{E} , which is given by:

$$\mathcal{E} = \left\{ |\hat{\mu}_t(a) - \mu(a)| \leq \sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_t(a)}}, \quad \text{for all } a = 1, 2, \dots, K \text{ and } t = 1, 2, \dots, T \right\}. \quad (6.6)$$

By Hoeffding's inequality, we know that with probability at least $1 - \delta$, we have

$$|\hat{\mu}_t(a) - \mu(a)| \leq \sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_t(a)}} \quad \text{for all } a = 1, 2, \dots, K \text{ and } t = 1, 2, \dots, T,$$

so that \mathcal{E} holds with probability at least $1 - \delta$. For an arm a , consider the total number of times it is played across T rounds, which can be described as $n_T(a)$ in our notation. Clearly, the regret of UCB is at most $\sum_{a=2, \dots, K} \Delta_a n_T(a)$. If T_a is the last round where a was played, then a maximized the UCB objective (6.5) on this round, which means in particular that

$$\hat{\mu}_{T_a}(a) + \sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_{T_a}(a) - 1}} \geq \hat{\mu}_{T_a}(1) + \sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_{T_a}(1)}},$$

where the inequality uses $n_{T_a}(a) = n_T(a) - 1$. Since \mathcal{E} holds, we further have

$$\mu(a) + 2\sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_T(a) - 1}} \geq \mu(1)$$

Simplifying, we see that

$$2\sqrt{\frac{\ln \frac{2KT}{\delta}}{2n_T(a) - 1}} \geq \Delta_a,$$

which gives the statement of the theorem with some algebra. ■

Note that the quantities Δ_a are not necessary for the algorithm, but only for the analysis. This is an instance of a problem dependent bound. If we want a bound purely in terms of K and T , then we can use the critical value of $\Delta_a = 1/\sqrt{KT}$. For each arm where Δ_a is larger than this value, then using a smaller value only yields a larger bound. If Δ_a is smaller than this value, we can use the worst possible upper bound of T for such an action, and the overall regret from this action is still at most $\sqrt{T/K}$. This line of reasoning says that the regret of UCB is always at most $\tilde{O}(\sqrt{KT})$, which is useful to know since our earlier bound seems to go to ∞ if Δ_a is very small.

6.4 Simple Strategies for i.i.d. Contextual Bandits

In order to gain further intuition about possible approaches to the contextual bandit problem, let us start from a simpler i.i.d. version. In this case, there is a fixed, but unknown joint distribution over contexts and rewards. We will use D to denote this joint distribution, as well as the marginal distribution of contexts. We further simplify the problem to a *full-information* version, where the player's interaction happens as described in Figure 6.3.

- At each round $t = 1, 2, \dots, T$:
1. Environment reveals a context $x_t \in \mathcal{X}$ sampled from the distribution D .
 2. Learner chooses an action $a_t \in \mathcal{A}$.
 3. Learner gains reward $r_t(a_t)$, and additionally observes $r_t(a)$ for all $a \in \mathcal{A}$, where $\mathbf{r} \sim D(\cdot|a, x)$.

Figure 6.3: The full-information version of i.i.d. contextual bandit problem

6.4.0.1 The GREEDY algorithm using full-information

The player's task in the full-information is considerably simpler, as there is no need to explore. While they still only gain the reward for the chosen action, the rewards of all the other actions are also revealed. This allows the algorithm to assess what might have happened if it chose a different action. Let us define the following GREEDY strategy:

1. At time t , let π_t be any policy which satisfies:

$$\pi_t = \arg \max_{\pi \in \Pi} \sum_{s=1}^{t-1} r_s(\pi(x_s)) \quad (6.7)$$

2. Predict $\pi_t(x_t)$.

We begin by showing why this strategy obtains a low regret for the i.i.d. setting. In order to prove the result, we will use the following form of Hoeffding's inequality, which can be obtained from the one you saw in the first lecture with some algebra.

Lemma 5 (Hoeffding's inequality). *Suppose X_1, \dots, X_n are independent and $a_i \leq X_i \leq b_i$. Let $\bar{X} = \sum_{i=1}^n X_i/n$. Then for any $\delta \in (0, 1]$, we have*

$$|\bar{X} - \mathbb{E}[\bar{X}]| \leq \sqrt{\frac{\sum_{i=1}^n (b_i - a_i)^2}{2n^2} \log \frac{2}{\delta}},$$

with probability at least $1 - \delta$. In the special case where $b_i \equiv 1$ and $a_i \equiv 0$ for all i , the following simpler bound holds with probability at least $1 - \delta$,

$$|\bar{X} - \mathbb{E}[\bar{X}]| \leq \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}.$$

Using the inequality, we now analyze the *instantaneous regret* of the GREEDY strategy first. The overall regret bound follows as a simple corollary.

Theorem 6. *The policy π_t computed at round t of the GREEDY algorithm satisfies, with probability at least $1 - \delta$*

$$V(\pi_t) \geq \max_{\pi \in \Pi} V(\pi) - \sqrt{\frac{2}{t} \log \frac{2N}{\delta}}.$$

Proof: For a fixed policy π , the random variables $r_s(\pi(x_s))$ for $s = 1, \dots, t$ are i.i.d. with mean $V(\pi)$, due to the i.i.d. assumption on the contexts and actions. Since the rewards are further assumed to be in $[0, 1]$ (recall Figure 6.1), we can invoke the simple form of Hoeffding's inequality from Lemma 5 to conclude that with probability at least $1 - \delta$,

$$\left| \frac{1}{t} \sum_{s=1}^t r_s(\pi(x_s)) - V(\pi) \right| \leq \sqrt{\frac{1}{2t} \log \frac{2}{\delta}}. \quad (6.8)$$

Unfortunately, we cannot apply the inequality directly to π_t , which is a policy that is picked after looking at the contexts and rewards over these t rounds. So the random variables $r_s(\pi_t(x_s))$ are not i.i.d. with mean $V(\pi_t)$. A standard trick to circumvent this difficulty is to establish a *uniform deviation* bound. That is, we show a bound similar to the one above, but simultaneously for all $\pi \in \Pi$. The idea is that if the sample averages and expectations are close for each $\pi \in \Pi$, then they will also be close for π_t in particular. To do this, let $\mathcal{E}(\pi)$ refer to the event that the bound in Equation 6.8 holds, and let $\mathcal{E}(\pi)^C$ be the complement of that event. Equation 6.8 then implies that for any fixed $\pi \in \Pi$, $\mathbb{P}(\mathcal{E}(\pi)^C) \leq \delta$. Observe that we are interested in the probability of $\mathcal{E}(\pi)$ holding for all $\pi \in \Pi$ simultaneously, which means that we want to control

$$\mathbb{P}(\cup_{\pi \in \Pi} \mathcal{E}(\pi)^C) \leq \sum_{\pi \in \Pi} \mathbb{P}(\mathcal{E}(\pi)^C) \leq N\delta.$$

That is, with probability $1 - N\delta$, Equation 6.8 holds simultaneously for all π , and for π_t in particular. Now, let π^* refer to the best policy in hindsight, that is $\pi^* = \arg \max_{\pi \in \Pi} V(\pi)$.¹ Then we will invoke Equation 6.8 twice, once with π_t and once with π^* . That is, with probability $1 - N\delta$, we have

$$\begin{aligned} \frac{1}{t} \sum_{s=1}^t r_s(\pi_t(x_s)) &\geq V(\pi_t) - \sqrt{\frac{1}{2t} \log \frac{2}{\delta}}. \\ V(\pi^*) &\geq \frac{1}{t} \sum_{s=1}^t r_s(\pi^*(x_s)) - \sqrt{\frac{1}{2t} \log \frac{2}{\delta}}. \end{aligned}$$

Adding the two inequalities, and noticing that the empirical reward of π_t is larger than that of π^* by definition (recall Equation 6.7), we obtain that with probability at least $1 - N\delta$

$$V(\pi_t) \geq \max_{\pi \in \Pi} V(\pi) - \sqrt{\frac{2}{t} \log \frac{2}{\delta}}.$$

Redefining δ/N as δ completes the proof. ■

As mentioned earlier, the regret bound for the *greedy algorithm* follows as a corollary.

Corollary 7. *With probability at least $1 - \delta$, the regret of the GREEDY algorithm is at most $\sqrt{8T \log \frac{2NT}{\delta}}$.*

The corollary follows by noting that the regret at round t is simply the suboptimality of the policy π_t which is controlled in Theorem 6. We can add the terms across rounds, and use the fact that $\sum_{t=1}^T 1/\sqrt{t} \leq 2\sqrt{T}$ to obtain the corollary.

We will show in the subsequent lectures how this simple strategy achieves low-regret in the full-information setting. We will also describe how this algorithm can be implemented efficiently, under a particular notion of efficiency.

¹Note that we can pick any policy that attains this largest value, if it is not unique.

6.4.0.2 The τ -GREEDY algorithm for bandit feedback

The simplicity of the GREEDY algorithm provokes a natural question. If something so simple works in the full-information setting, is there a natural counterpart in the partial information setting where we only observe the reward of the chosen action. We now describe a relatively direct extension of the GREEDY strategy, called τ -GREEDY which can be executed in a contextual bandit setting. We first introduce the notation $\mathbf{1}(\cdot)$ to be the 0 – 1 indicator for whether the argument of it is true or not. The τ -greedy algorithm for contextual bandits is presented in Algorithm 1.

Algorithm 1 The τ -GREEDY algorithm for contextual bandits

Require: Exploration parameter $\tau \in (0, T]$.

Initialize $S = \emptyset$ to be the empty set.

for $t = 1, 2, \dots, \tau$ **do**

 Play a_t uniformly at random from $1, \dots, K$.

 Update $S = S \cup \{(x_t, a_t, r_t)\}$.

end for

Let π_τ be any policy which satisfies:

$$\pi_\tau = \arg \max_{\pi \in \Pi} \sum_{(x, a, r) \in S} r \mathbf{1}(\pi(x) = a). \quad (6.9)$$

for $t = \tau + 1, \dots, T$ **do**

 Play $a_t = \pi_\tau(x_t)$.

end for

The τ -GREEDY strategy is effectively trying to mimic the GREEDY strategy. It deals with the partial feedback by doing uniformly random exploration on an initial subset of the rounds, and uses the data collected on these rounds to assess the quality of each policy. The greedy policy based on this data is then used for the remaining rounds, leading to a clear demarcation between the *explore* and *exploit* phase.

We will now give a regret bound for this algorithm, which formalizes the above intuition.

Theorem 8. *The regret of the τ -GREEDY algorithm shown in Algorithm 1 is bounded with probability at least $1 - \delta$ as $\mathcal{O}((TK)^{2/3}(\log(N/\delta))^{1/3})$.*

Proof: Note that the policy π_τ computed in Equation 6.9 is similar to the policy computed in Equation 6.7, except that we only compute it once and on τ examples only. We begin with a slight rewriting of the maximization criterion in Equation 6.9. Let us define the objective at time t as

$$\widehat{R}_\tau(\pi) = \sum_{s=1}^{\tau} r_s \mathbf{1}(\pi(x_s) = a_s).$$

Taking an individual summand, let us define $Y_s = r_s \mathbf{1}(\pi(x_s) = a_s)$. Then we can compute its expectation as

$$\begin{aligned} \mathbb{E}[Y_s] &= \mathbb{E}[r_s \mathbf{1}(\pi(x_s) = a_s)] \\ &= \mathbb{E}[\mathbb{E}[r_s \mathbf{1}(\pi(x_s) = a_s) \mid x_s]] \\ &\stackrel{(a)}{=} \frac{1}{K} \mathbb{E}[\mathbb{E}[r(\pi(x_s)) \mid x_s]] \\ &= \frac{1}{K} V(\pi). \end{aligned}$$

Here (a) follows since the identity of the chosen action in the explore events is independent of x_s as well as r_s . Hence, adding up the expectations, we see that

$$\mathbb{E}[\widehat{R}_\tau(\pi)] = \frac{\tau}{K} V(\pi),$$

and that it is indeed a sum of τ i.i.d. terms with this mean. That is, we are effectively working with τ/K samples instead of τ samples in defining the greedy policy. We can apply Hoeffding's inequality as in the proof of Theorem 6 to $\widehat{R}_t(\pi)$ now, as well as combine it with union bound, to see that with probability at least $1 - \delta$, simultaneously for all $\pi \in \Pi$

$$\left| \frac{1}{\tau} \widehat{R}_\tau(\pi) - \frac{V(\pi)}{K} \right| \leq \sqrt{\frac{1}{2\tau} \log \frac{2N}{\delta}}.$$

Now we apply the inequality twice as before, once with π_τ and once with π^* to conclude that

$$V(\pi_\tau) \geq V(\pi^*) - K \sqrt{\frac{2}{\tau} \log \frac{2N}{\delta}}.$$

Note that the regret of the algorithm is now written as the suboptimality of $\pi(\tau)$ on rounds $\tau + 1, \dots, T$ and can be arbitrarily large on the first τ rounds. Hence we bound the regret by

$$\tau + (T - \tau)(V(\pi_\tau) - V(\pi^*)) \leq \tau + TK \sqrt{\frac{2}{\tau} \log \frac{2N}{\delta}}.$$

Minimizing over τ now yields the theorem statement. ■

This regret bound is significantly inferior than that of UCB in the non-contextual setting, in that it scales as $\mathcal{O}(T^{2/3})$ instead of $\mathcal{O}(\sqrt{T})$. Can we develop better algorithms which are also computationally efficient in the i.i.d. setting? The next lecture will take a slight digression and introduce some necessary concepts to help us answer this question. We will also see how algorithms like GREEDY and τ -GREEDY can be implemented in a computationally efficient manner in the next lecture.

References