

Advanced Topics in Machine Learning 2017-2018

Yevgeny Seldin Christian Igel Tobias Sommer Thune Julian Zimmert

Home Assignment 7

Deadline: Tuesday, 31 October 2017, 23:59

The assignments must be answered individually - each student must write and submit his/her own solution. We encourage you to work on the assignments on your own, but we do not prevent you from discussing the questions in small groups. If you do so, you are requested to list the group partners in your individual submission.

Submission format: Please, upload your answers in a single .pdf file and additional .zip file with all the code that you used to solve the assignment. (The .pdf should **not** be part of the .zip file.)

IMPORTANT: We are interested in how you solve the problems, not in the final answers. Please, write down the calculations and comment your solutions.

1 Label-efficient prediction and label-efficient bandits (33 points)

In many applications getting labels is expensive. In this question we investigate what can be done in adversarial settings when we are restricted by the number of labels we are allowed to request.

We start with adversarial full information setting. Assume that the game lasts for T rounds and the algorithm is allowed to observe on average εT columns in the matrix of losses for $0 < \varepsilon \leq 1$ (which means that it can evaluate the quality of the experts it is working with on εT out of T rounds). Consider the playing strategy in Algorithm 1:

Algorithm 1 Label-Efficient Forecaster

```
 $\forall a : \tilde{L}_0(a) = 0$ 
for  $t = 1, \dots, T$  do
   $\forall a : p_t(a) = \frac{e^{-\eta \tilde{L}_{t-1}(a)}}{\sum_{a'} e^{-\eta \tilde{L}_{t-1}(a' )}}$ 
  Sample  $A_t$  according to  $p_t$  and play it
  Draw Bernoulli random variable  $Z_t$  with bias  $\varepsilon$ 
  if  $Z_t = 1$  then
    Request to observe  $\ell_t^1, \dots, \ell_t^K$ 
  else
    Make no observations
  end if
   $\forall a : \tilde{\ell}_t^a = \frac{\ell_t^a \mathbb{1}\{Z_t=1\}}{\varepsilon} = \begin{cases} \frac{\ell_t^a}{\varepsilon}, & \text{if } Z_t = 1 \\ 0, & \text{otherwise} \end{cases}$ 
   $\forall a : \tilde{L}_t(a) = \tilde{L}_{t-1}(a) + \tilde{\ell}_t^a$ 
end for
```

1. Show that in expectation the above algorithm requests to observe εT columns.
2. Show that the expected regret of the above algorithm in adversarial setting satisfies $\mathbb{E}[R_T] \leq \sqrt{2\frac{1}{\varepsilon} T \ln K}$. Note that for $\varepsilon = 1$ it is allowed to observe all columns and we recover the Hedge algorithm, whereas for $\varepsilon < 1$ the performance gradually degrades with $\sqrt{\frac{1}{\varepsilon}}$.

3. Show that with high probability the exact number of observed columns is not significantly larger than εT .
4. How would you modify the algorithm if you had a hard restriction εT on the number of columns that you are allowed to observe? How would you analyze the modified algorithm? (*Hint: you know that if you have run over the budget before the game finished your regret is still bounded by T . So all you have to do is to make sure that you do not run over the budget too often. The exact parameter tuning in this case is a bit nasty, so you are allowed to give an answer in a parametrized form.*)
5. Now consider an adversarial *bandit* setting and assume that an algorithm allowed to observe its own loss on average on εT out of T rounds of the game. Propose an algorithm for this setting and analyze its expected regret.

You are welcome to evaluate the algorithms empirically (not for submission). I.i.d. experiments will be the easiest to construct, but you are also welcome to try adversarial settings.

2 Policy evaluation (33 points)

Proof that *iterative policy evaluation* with update

$$\forall s \in S : V_{k+1}(s) \leftarrow \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V_k(s')]$$

in step k and $0 < \gamma < 1$ converges to the true value function, $\forall s \in S : \lim_{k \rightarrow \infty} V_k(s) = V^\pi(s)$. This can be done using the tools we have worked with during the course. One way is to consider

$$\Delta_k = \max_s |V_k(s) - V^\pi(s)| .$$

3 Reinforcement learning (34 points)

Let us consider the navigation task defined by depicted floor plan in Figure 1. The floorplan represents a 3×4 grid-world. Each of the 12 rooms corresponds to one state. The agent can go from one room to another if the rooms are connected by a door. The actions are the movements {up, down, right, left}. Every movement (except in the terminal state) gives a negative reward of -1 (i.e., every movements costs 1). If you bump into a wall without a door, you stay in the same room, but the movement has still to be paid. There are three rooms that give you an *additional* negative reward of -5 and -10 , respectively, *when you enter them*, see figure. The room in the bottom right is a terminal state. An episode ends when this room is reached, which can be modelled by every action in this state leaving the state unchanged and having no cost.

This exercise requires an implementation of the MDP. Note that all rewards and transitions are deterministic. They could be described by simple mappings $S \times A \rightarrow \mathbb{R}$ and $S \times A \rightarrow S$, respectively, which can be encoded by simple tables.

1. Use an algorithm presented in the lecture to compute the value function V^{rand} of the random policy, that is, the policy that chooses an action uniformly at random in every state. Report the 12 V^{rand} values. Provide the implementation of the algorithm.
2. Use an algorithm presented in the lecture to compute the optimal value function V^* . Report the 12 V^* values. Provide the implementation of the algorithm.

4 Bonus: Make Your Own Question

Propose a question on any of the course material. If you manage to come up with a good question, it may be used in future editions of the course and you may get some bonus points.

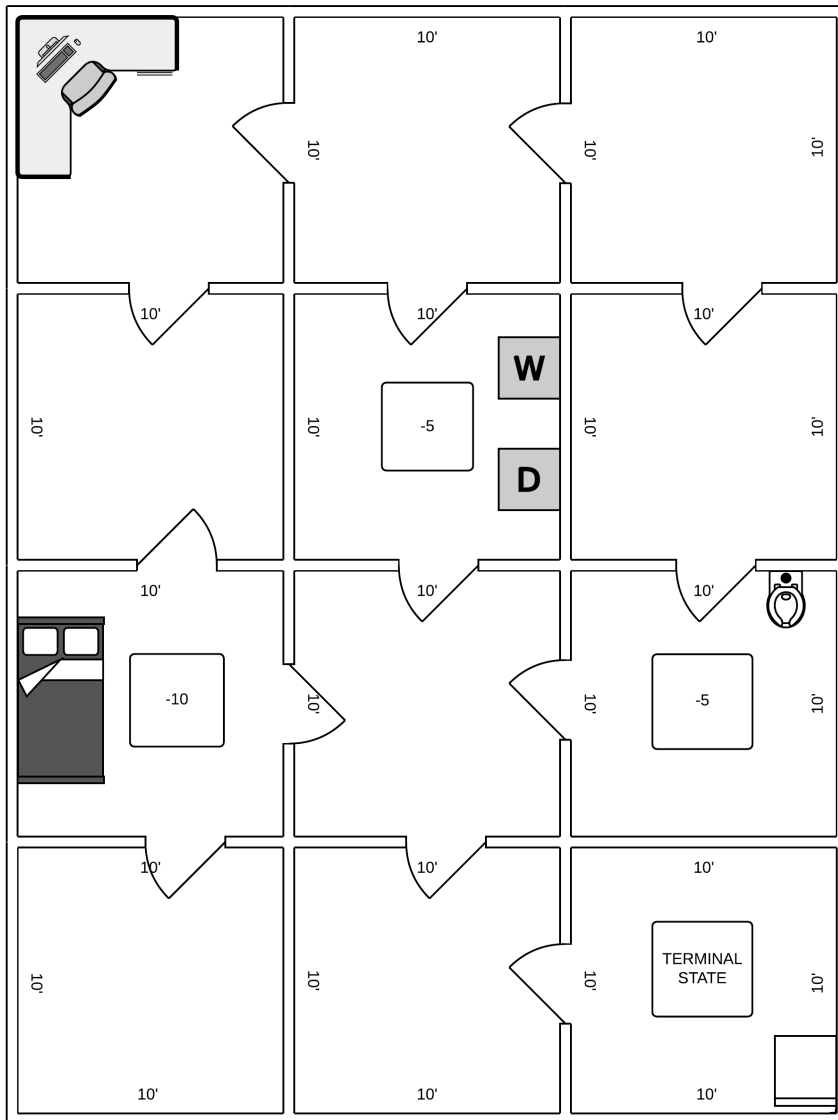


Figure 1: Floor plan.

5 [Optional, not for submission] Rewards vs. losses

The original EXP3 algorithm for multiarmed bandits was designed for the game with rewards rather than losses (see Auer et al. (2002, Page 6)). In the game with rewards we have an infinite matrix of rewards $\{r_t^a\}_{a \in \{1, \dots, K\}, t \geq 1}$, where $r_t^a \in [0, 1]$. On each round of the game the algorithm plays an action A_t and accumulates and observed reward $r_t^{A_t}$. The remaining rewards r_t^a for $a \neq A_t$ remain unobserved.

On the one hand, we can easily convert rewards into losses by taking $\ell_t^a = 1 - r_t^a$ and apply the EXP3 algorithm we saw in class. On the other hand, a bit surprisingly, working directly with rewards (as Auer et al. did) turns to be more cumbersome and less efficient. The high-level reason is that the games with rewards and losses have a different dynamics. In the rewards game when an action is played its relative quality (expressed by the cumulative reward) increases. Therefore, we need explicit exploration to make sure that we do not get locked on a suboptimal action. In the losses game when an action is played its relative quality (expressed by the cumulative loss) decreases. Therefore, we never get locked on any particular action and exploration happens automatically without the need to add it explicitly (sometimes this is called implicit exploration). The low-level reason when it comes down to the analysis of the algorithm is that it is easier to upper bound the exponent of x for negative x as opposed to positive x .

The original EXP3 algorithm for the rewards game looks as follows, where the most important difference with the algorithm for the losses game is highlighted in red and two additional minor differences are highlighted in blue (we explicitly emphasize that the sign in the exponent changes from “-” to “+”). \tilde{R}_t is used to denote cumulative importance-weighted rewards.

Algorithm 2 The EXP3 Algorithm for the game with rewards and fixed time horizon

```

1:  $\forall a : \tilde{R}_0(a) = 0$ 
2: for  $t = 1, \dots, T$  do
3:    $\forall a : p_t(a) = (1 - \eta) \frac{e^{+\eta \tilde{R}_{t-1}(a)}}{\sum_{a'} e^{+\eta \tilde{R}_{t-1}(a')}} + \frac{\eta}{K}$ 
4:   Sample  $A_t$  according to  $p_t$  and play it
5:   Observe  $r_t^{A_t}$ 
6:    $\forall a : \tilde{r}_t^a = \frac{r_t^a \mathbb{1}\{A_t=a\}}{p_t(a)} = \begin{cases} \frac{r_t^a}{p_t(a)}, & \text{if } A_t = a \\ 0, & \text{otherwise} \end{cases}$ 
7:    $\forall a : \tilde{R}_t(a) = \tilde{R}_{t-1}(a) + \tilde{r}_t^a$ 
8: end for
```

1. Explain why the analysis of the EXP3 algorithm for the rewards game without the addition of explicit exploration $\frac{\eta}{K}$ in Line 3 of the algorithm would not work. More specifically - if you would try to follow the lines of the analysis of EXP3 with losses, at which specific point you would get stuck and why?
2. How the addition of explicit exploration term $\frac{\eta}{K}$ in Line 3 of the algorithm allows the analysis to go through? (You can check the analysis of the algorithm in Auer et al. (2002, Page 7).)
3. By how much the expected regret guarantee for EXP3 with rewards is weaker than the expected regret guarantee for EXP3 with losses? (Check Auer et al. (2002, Corollary 3.2) and assume that g takes its worst-case value, which is T .)
4. You are mostly welcome to experiment and see whether theoretical analysis reflects the performance in practice. I.i.d. experiments will be the easiest to construct, but you are also welcome to try adversarial settings.

Good luck!
Yevgeny, Christian, Tobias & Julian

References

Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal of Computing*, 32(1), 2002.