

Homework 2

6501 Reinforcement Learning (Spring 2025)

Submission deadline: 11:59pm, February 26

Contextual Bandit Algorithms

In this homework, we will implement bandit algorithms we have covered in the lecture. Specifically, we will implement

- Algorithms based on regression oracle and exploration mechanisms including ϵ -greedy (EG), Boltzmann exploration (BE), and inverse gap weighting (IGW), and
- Proximal Policy Optimization (PPO).

The starter code can be accessed at http://bahh723.github.io/rl2025sp_files/HW2.py.

The latex template can be accessed at <https://www.overleaf.com/read/bktngfjrszmx#3e0687>

Submit the answer as .pdf and the code as a .py file to Gradescope.

1 Data

To simulate a contextual bandit environment, we use existing classification dataset for supervised learning. Specifically, in this homework, we use the [mnist dataset](#) with some pruning and modification. Originally, it is a classification dataset where features are $28 \text{ pixel} \times 28 \text{ pixel}$ gray-scale images of digits, and the classes correspond to 10 digits.

2 Data Conversion

For simplicity, we trim the dataset so that it only contains 4 classes corresponding to digits ‘0’, ‘1’, ‘2’, and ‘3’, each having 2500 samples, summing up to 10000 samples. In each round, the environment will reveal the image (context), and the learner has to pick one of the classes (action). We artificially make the reward function change once in the middle. Specifically, in our case the time horizon is $T = 10000$. For $t \leq 5000$ (Phase 1), the reward function is the following:

$$R(x, a) = \begin{cases} 0.5 & \text{if } a \text{ is the correct digit of image } x \\ 0 & \text{otherwise} \end{cases}$$

For $t > 5000$ (Phase 2), the reward is the following:

$$R(x, a) = \begin{cases} 0.5 & \text{if } a \text{ is the correct digit of image } x \\ 1 & \text{if } (a - 1) \bmod 4 \text{ is the correct digit of image } x \\ 0 & \text{otherwise} \end{cases}$$

For example, in Phase 2, if the learner chooses action 3 when seeing an image of digit 2, then it receives a reward of 1. We assume that the learner sees the true noiseless reward $R(x_t, a_t)$.

This conversion is already implemented in the starter code, so you don’t need to do anything here.

3 Algorithms based on regression oracle

In this part, we will implement the simple value-based algorithm based on regression on the reward function. In the class, we have introduced three common exploration mechanisms in this case: ϵ -Greedy, Boltzmann Exploration, and Inverse Gap Weighting. They share the same pseudo-code outlined in [Algorithm 1](#).

Algorithm 1 Algorithms based on regression

```

1 Randomly initialize a reward network  $R_\theta$  that takes the context as input and outputs the reward of each action.
2 Let  $\theta_1$  be the initial weights for the reward network.
3 for  $t = 1, \dots, T$  do
4   for  $n = 1, \dots, N$  do
5     Receive context  $x_{t,n}$ .
6     Sample action  $a_{t,n} \sim \pi(\cdot \mid x_{t,n})$  where  $\pi(\cdot \mid x) = f(R_{\theta_t}(x, \cdot))$ 
7     Receive reward  $r_{t,n}$ .
8    $\theta \leftarrow \theta_t$ 
9   for  $m = 1, \dots, M$  do
10    
$$\theta \leftarrow \theta - \lambda \nabla_\theta \left( \frac{1}{N} \sum_{n=1}^N (R_\theta(x_{t,n}, a_{t,n}) - r_{t,n})^2 \right). \quad (1)$$

11   $\theta_{t+1} \leftarrow \theta$ 

```

In [Algorithm 1](#), $f: \mathbb{R}^A \rightarrow \Delta_A$ is a link function that maps an A -dimensional real vector to a distribution over A actions. The three exploration mechanisms correspond to the following three choices of f :

- ϵ -Greedy: $[f(v)]_a = \frac{\epsilon}{A} + (1 - \epsilon)\mathbb{I}\{a = \operatorname{argmax}_{a'} v(a')\}$.
- Boltzmann Exploration: $[f(v)]_a = \frac{e^{\lambda v(a)}}{\sum_{a'} e^{\lambda v(a')}}.$
- Inverse Gap Weighting: $[f(v)]_a = \begin{cases} \frac{1}{A + \lambda(\max_{a'} v(a') - v(a))} & \text{if } a \neq \operatorname{argmax}_{a'} v(a') \\ 1 - \sum_{a' \neq a} [f(v)]_{a'} & \text{if } a = \operatorname{argmax}_{a'} v(a') \end{cases}.$

[Algorithm 1](#) may be slightly different from the standard way of implementing such algorithm — usually, in [Eq. \(1\)](#) we may reuse data collected in the past by having a *replay buffer* that stores old data. But since our artificial problem has a non-stationary reward function, reusing data from the past becomes not a great idea.

Please complete the following tasks in (a)-(g). Note that the starter code already calculates the average reward in Phase 1, Phase 2, and overall horizon, respectively. The starter code also already implements the regression oracle (1). A figure of running average is also generated for a single parameter. The major task is to code up the three link functions above, run the experiments, and combine the figures for multiple parameters.

Note that you are allowed to change the default hyperparameters in the starter code (N , M , optimizer, learning rates, etc.). In the table below, you may also change the values of hyperparameters or add additional ones if you feel that the given values cannot reflect the trend.

- (a) (5%) Implement ϵ -Greedy and, for different values of ϵ , record in the table below the average reward in Phase 1, Phase 2, and over the entire horizon.

ϵ	Phase 1	Phase 2	Overall
0.1			
0.03			
0.01			
0.003			
0.001			
0			

- (b) (5%) Plot the running average reward over time in the same figure for your experiments in (a). See appendix for an example of such a figure.
- (c) (5%) Implement Boltzmann Exploration and, for different values of λ , record in the table below the average reward in Phase 1, Phase 2, and over the entire horizon.

λ	Phase 1	Phase 2	Overall
2			
5			
10			
20			
50			

- (d) (5%) Plot the running average reward over time in the same figure for your experiments in (c).
- (e) (5%) Implement Inverse Gap Weighting and, for different values of λ , record in the table below the average reward in Phase 1, Phase 2, and over the entire horizon.

λ	Phase 1	Phase 2	Overall
10			
30			
100			
300			
1000			

- (f) (5%) Plot the running average reward over time in the same figure for your experiments in (e).

- (g) (5%) Have you noticed any trends in the experiments from (a)–(f)? In particular, how does the parameter influence the average reward in Phase 1 and Phase 2, respectively? What is the underlying explanation for this effect?

4 PPO

Algorithm 2 PPO without clipping

```

11 Randomly initialize a policy network  $\pi_\theta$  that takes contexts as input and outputs an action distribution.
12 Let  $\theta_1$  be the initial weights for the policy network.
13 If using adaptive baseline, additionally initialize a baseline network  $b_\phi$ .
14 for  $t = 1, \dots, T$  do
15   for  $n = 1, \dots, N$  do
16     Receive context  $x_{t,n}$ .
17     Sample action  $a_{t,n} \sim \pi_{\theta_t}(\cdot | x_{t,n})$ 
18     Receive reward  $r_{t,n}$ .
19    $\theta \leftarrow \theta_t$ 
20   for  $m = 1, \dots, M$  do
21     
$$\theta \leftarrow \theta + \lambda \nabla_\theta \left\{ \frac{1}{N} \sum_{n=1}^N \left[ \frac{\pi_\theta(a_{t,n} | x_{t,n})}{\pi_{\theta_t}(a_{t,n} | x_{t,n})} (r_{t,n} - b_{t,n}) - \frac{1}{\eta} \left( \frac{\pi_\theta(a_{t,n} | x_{t,n})}{\pi_{\theta_t}(a_{t,n} | x_{t,n})} - 1 - \log \frac{\pi_\theta(a_{t,n} | x_{t,n})}{\pi_{\theta_t}(a_{t,n} | x_{t,n})} \right) \right] \right\}, \quad (2)$$

22     where
23     
$$b_{t,n} = \begin{cases} b & \text{for fixed baseline} \\ b_\phi(x_{t,n}) + b' & \text{for adaptive baseline} \end{cases} \quad (3)$$

24     If using adaptive baseline, update
25     
$$\phi \leftarrow \phi - \lambda' \nabla_\phi \left[ \frac{1}{N} \sum_{n=1}^N (b_\phi(x_{t,n}) - r_{t,n})^2 \right]. \quad (4)$$

26    $\theta_{t+1} \leftarrow \theta$ 

```

4.1 Fixed Baseline

- (a) (5%) Implement Algorithm 2 with fixed baseline (so Eq. (4) can be omitted) and, for different values of b , record in the table below the average reward in Phase 1, Phase 2, and over the entire horizon.

b	Phase 1	Phase 2	Overall
2			
1.5			
1			
0.5			
0			
-0.5			

- (b) (5%) Plot the running average reward over time in the same figure for your experiments in (a).
- (c) (5%) Have you noticed any trends in the experiments from (a)–(b)? How does the baseline affect the average reward in Phase 1 and Phase 2, respectively? What is the underlying explanation for this effect?

4.2 Adaptive Baseline

- (a) (5%) Implement [Algorithm 2](#) with adaptive baseline and, for different values of the extra baseline b' in (3), record in the table below the average reward in Phase 1, Phase 2, and over the entire horizon.

b'	Phase 1	Phase 2	Overall
0.2			
0.1			
0			
-0.1			
-0.2			

- (b) (5%) Plot the running average reward over time in the same figure for your experiments in (a).
- (c) (5%) Is there any difference between adaptive baseline and fixed baseline in [Section 4.1](#)? What are the potential advantages or disadvantages of using adaptive baseline?
- (d) (5%) What is the effect of the extra baseline parameter b' ?

5 Survey

- (a) (5%) How much time did you use to complete each of the problems in this homework? Do you have any suggestion for the course? (e.g., the pace of the lecture, the length of the homework)

A Appendix

In the starter code, we plot the learning curve for a single parameter. You need to plot the learning curves for multiple parameter values on the same figure. Below is an example that shows how to plot two curves, although in usually need to plot 5-6 curves in the problems above.

