

# Learning Dynamics: Assignment 3

## Multi-Armed Bandits and Stochastic Reward Game

Hakim Boulahya

hboulahy@ulb.ac.be - 000391737

Université Libre de Bruxelles

December 18, 2017

### Contents

<b>1</b>	<b>N-Armed Bandit</b>	<b>2</b>
1.1	Exercice 1 . . . . .	2
1.1.1	Average rewards . . . . .	2
1.1.2	Q-values estimation . . . . .	2
1.1.3	Histograms . . . . .	4
1.2	Exercice 2 . . . . .	7
1.2.1	Average rewards . . . . .	7
1.2.2	Q-values estimation . . . . .	7
1.2.3	Histograms . . . . .	9
1.3	Exercice 3 . . . . .	10
1.3.1	Average rewards . . . . .	10
1.3.2	Q-values estimation . . . . .	10
1.3.3	Histograms . . . . .	12
<b>2</b>	<b>Climbing game</b>	<b>12</b>

# 1 N-Armed Bandit

**Remark** About the notation, the first arm/action starts at #0.

## 1.1 Exercice 1

### 1.1.1 Average rewards

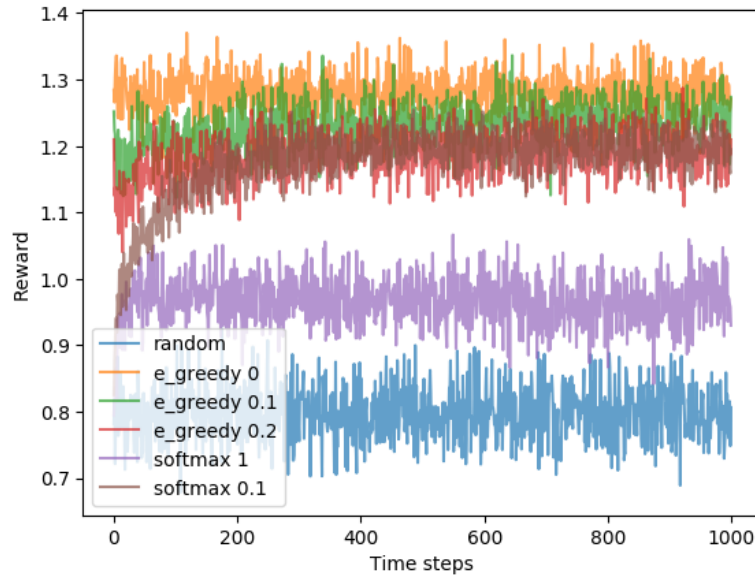


Figure 1: Average rewards for all algorithms

We can see that only exploring *i.e.* random and softmax with larger temperate  $\tau = 1$  doesn't perform well, and the reward stays low. Using softmax with a temperature of 0.1 gives better results because of the fact that it gets more greedy.

Using  $\epsilon$ -greedy with small value, *i.e.* a large probability to choose the optimal actions, tends to better rewards.  $\epsilon$ -greedy with 0 seems to arrives *directly* to the best arm. We can see on the graph that when  $\epsilon$  is bigger it usually takes more time to get to the best arm, because of the exploration. Softmax with a temperature of 0.1 does tends to the best arm, but takes more time.

It is important to note that the greedy behaviour is better, because the standard deviation is small for this exercice. Because of that, using the best action based on the estimated Q-values is accurate.

### 1.1.2 Q-values estimation

Figure 2 shows the estimation of the Q-values for all algorithms.

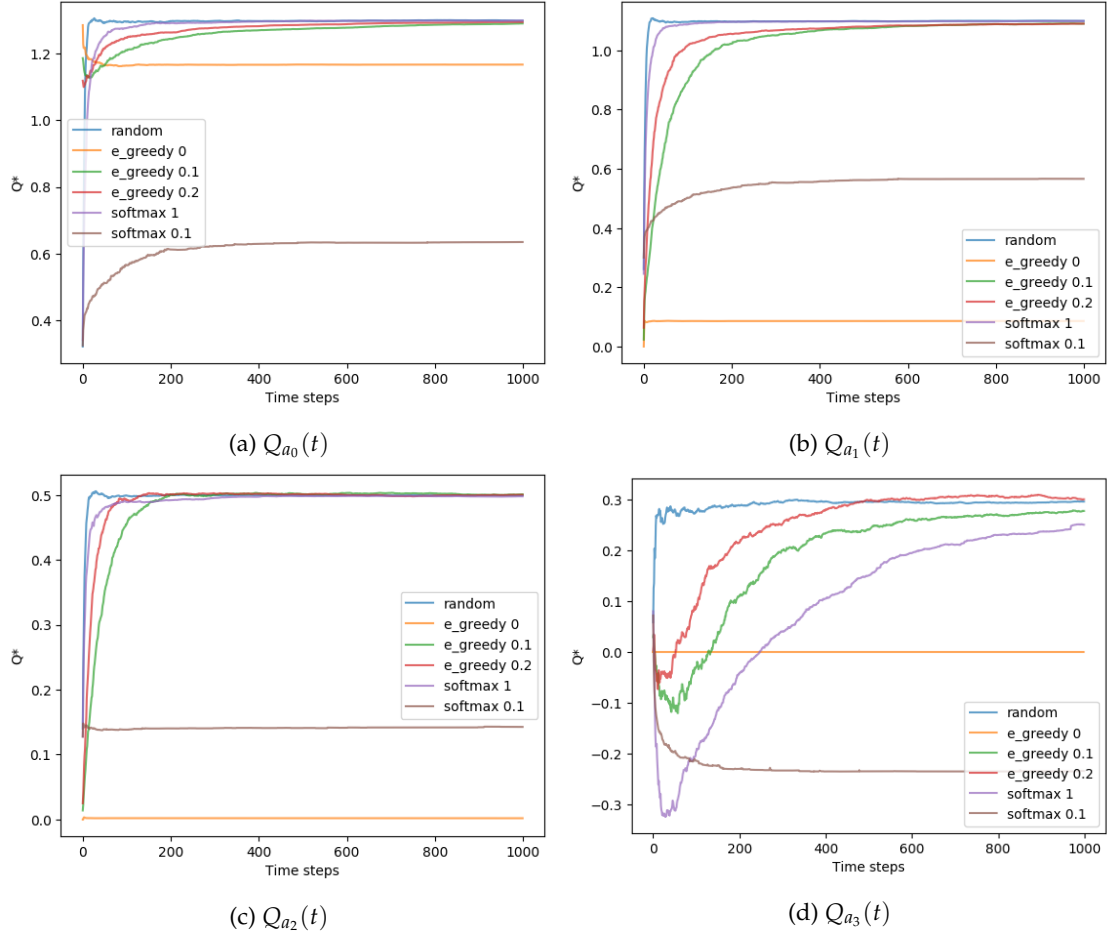


Figure 2: Plot per arm showing the  $Q_{a_i}^*$  of that action along with the actual  $Q_{a_i}$  estimate over time with  $\mu = (1.3, 1.1, 0.5, 0.3)$ ,  $\sigma = (0.9, 0.6, 0.4, 2.0)$

## Random

Since random methods is full exploration, we can see that for all Q-values the estimation of the  $Q_{a_i}^*$  are accurate.

## $\epsilon$ -greedy

$\epsilon = 0$  Using a *only* greedy exploration, we can see that only the best actions have Q-values different from 0. Actually the best action #0 tends to be estimated fast enough. For action #1, the Q-values is bigger than 0 because at the beginning of the learning, the algorithm is still looking for the greedy action, therefore action #1 could be the best one, until the algorithm find that #action 0 is better. The two other actions, are never explored since the it is a full greedy method.

$\epsilon = \{0.1, 0.2\}$  The estimation is more accurate when  $\epsilon > 0$  *i.e.* when the action are also selected randomly which allows exploration. We can see that by using a non-null  $\epsilon$ , the estimation of

the Q-values overtime tends to be accurate. An interesting observation is that the estimation of the Q-values is faster (the accuration rate) when  $\epsilon$  is larger. This is logical, since more  $\epsilon$  is large, more the action are selected randomly.

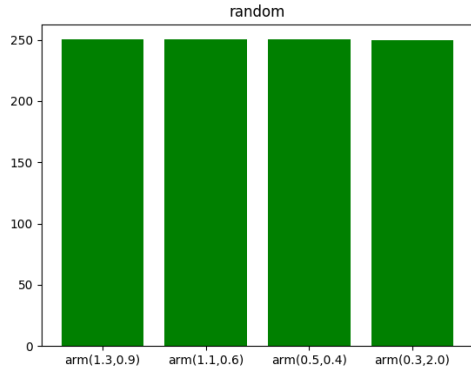
## Softmax

$\tau = 1$  Using  $\tau = 1$  (a large enough temperature) the Q-values estimation follows the same pattern as for any exploration methods like random or  $\epsilon = 0$ . We can see that except for the last action the curves are close to the random ones, which confirm the exploration nature of a large temperature. The fact that the curves of the last action #3 are more biased is because the standard deviation is larger than for the other actions, which make the Q-values estimation harder to make.

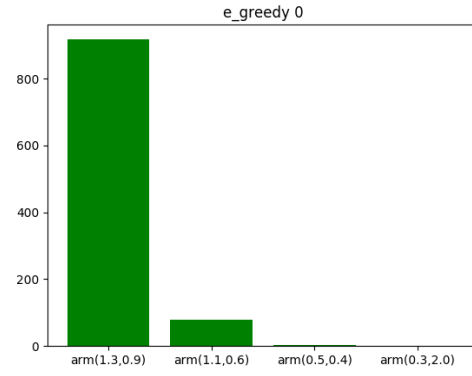
$\tau = 0.1$  Using a smaller temperature for the softmax selection method, we have a more greedy action selection. This action selection method is the worst one that estimate Q-values (excluding the full exploiting method 0-greedy). We can see that for the best arms #0 and #1, the estimations are larger than for there worst arms, but it is *far* from the actual  $Q^*$ . This can be explained by the nature of the softmax selection method with a small  $\tau$ : it sticks to the best arm found and doesn't explore much to estimate the correct Q-values overtime.

### 1.1.3 Histograms

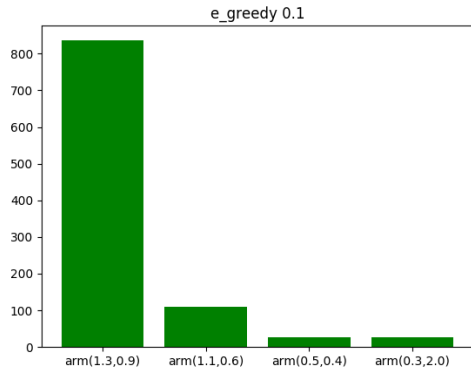
Figure 3 shows the histograms for all algorithms.



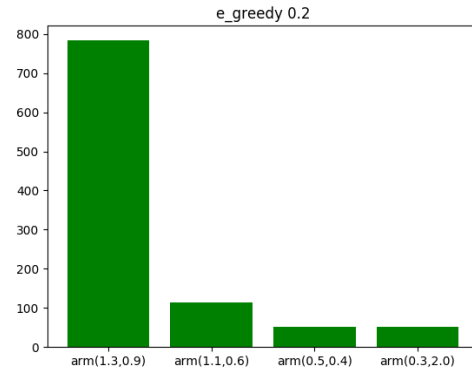
(a)



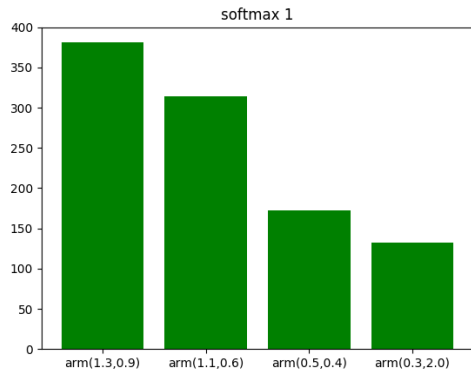
(b)



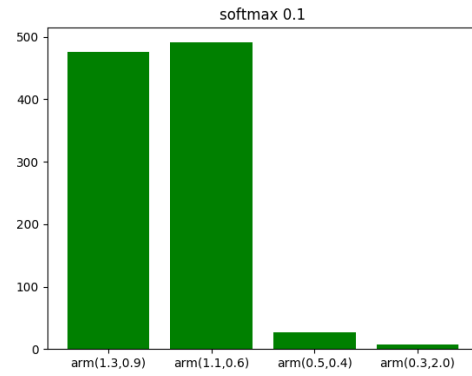
(c)



(d)



(e)



(f)

Figure 3: Histograms showing the number of times each action is selected per selection strategy with  $\mu = (1.3, 1.1, 0.5, 0.3)$ ,  $\sigma = (0.9, 0.6, 0.4, 2.0)$

## Random

The action selections, shown in Figure 3a, for the random selection method is equidistributed. Since all actions have the same probability it is logical that the selection is equidistributed when there  $t \rightarrow \infty$ .

## $\epsilon$ -greedy

Figures 3b 3c 3d shows the histograms of action selections for  $\epsilon$ -greedy selection methods.

$\epsilon = 0$  With a null  $\epsilon$ , the action selection will be greedy *i.e.*, and since we have a pretty small standard deviation, the best action #0 is always chosen when found by the algorithm. This means that there is no exploration, the method chooses the best arm directly. The fact that the other actions are chosen, is that at the beginning it is necessary to explore a little bit to find the best action.

$\epsilon = \{0.1, 0.2\}$  It follows the same pattern that when  $\epsilon$  is null, except that here it is possible for the methods to explore more *i.e.* choosing other actions than the best one. This is why the bar of other actions are higher when  $\epsilon$  is bigger. If  $\epsilon = 1$  the bars will be as shown for the random methods in Figure 3a.

## Softmax

Figures 3e 3f shows the histograms of action selections for softmax selection methods.

When temperature  $\tau = 1$ , the action selections are more randomize *i.e.* the exploration is more noticable. An interesting observation in Figure 3f is that when  $\tau = 0.1$ , *i.e.* the action selection is more greedy, softmax algorithm tends to hesitate to choose the best actions. The number of times that the 2 best actions, action #0 et #1, are chosen is more or less equal.

## 1.2 Exercice 2

### 1.2.1 Average rewards

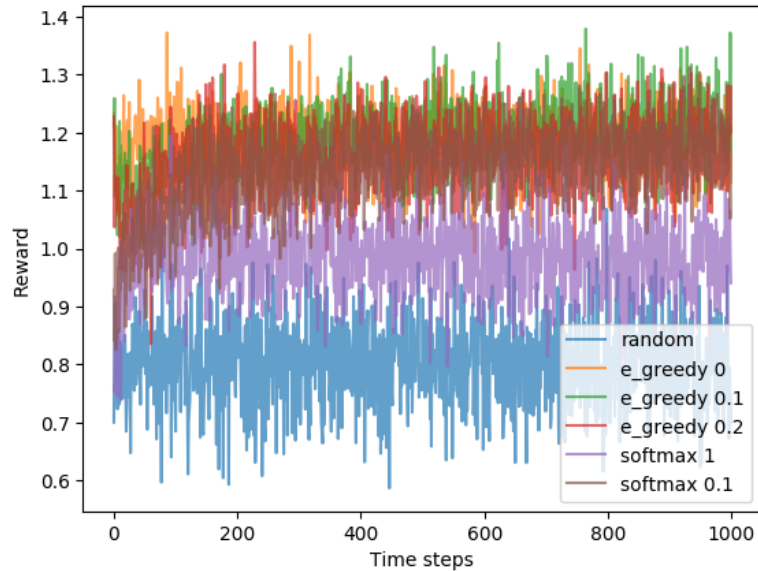


Figure 4: Average rewards for all algorithms

The results order of the algorithms that perform best doesn't really change compared to exercise 1. By doubling the standard deviation, we only allowed the algorithms to provide a more random rewards *i.e.* a bigger range. We can see that the dynamics of the algorithms don't change.

It is harder to learn because of the standard deviation that provides a more randomize outcomes, therefore harder to find the best rewards.

### 1.2.2 Q-values estimation

Figure 5 shows the estimation of the Q-values for all algorithms.

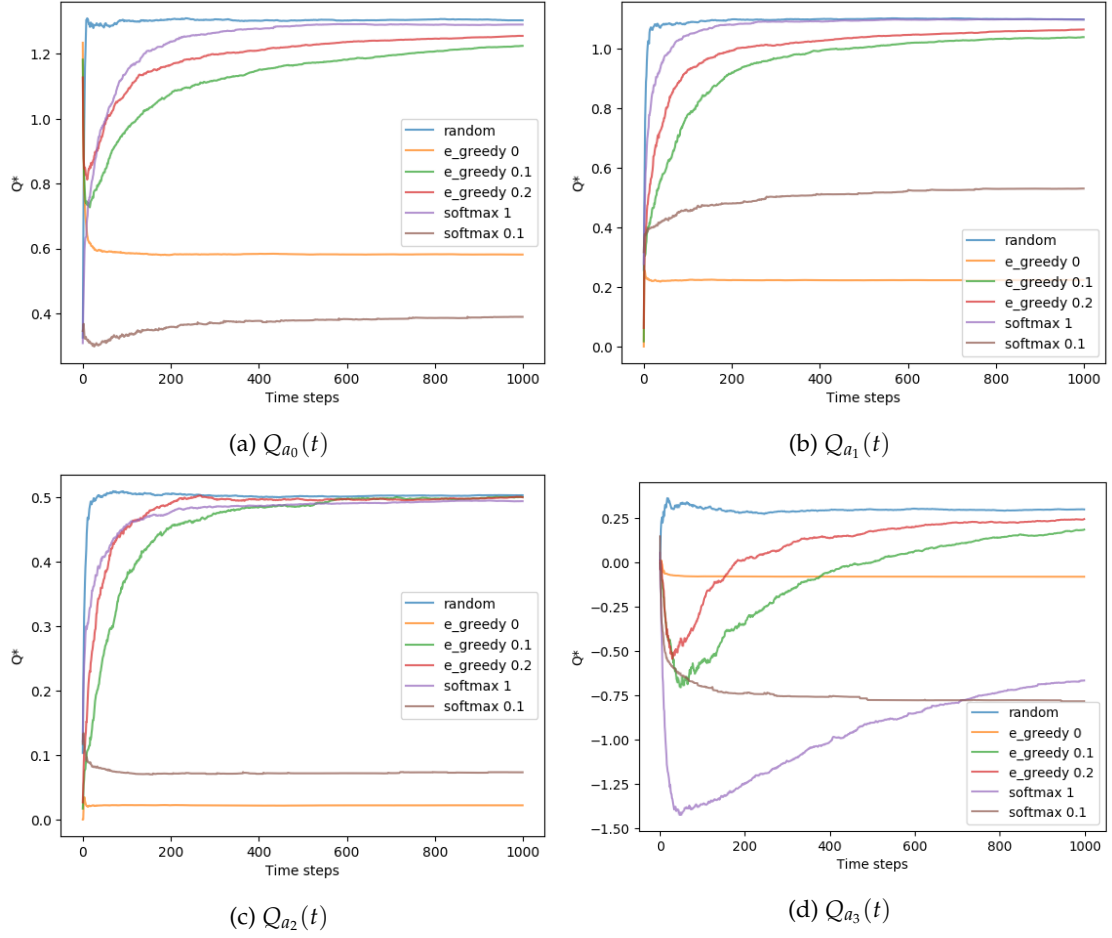


Figure 5: Plot per arm showing the  $Q_{a_i}^*$  of that action along with the actual  $Q_{a_i}$  a i estimate over time with  $\mu = (1.3, 1.1, 0.5, 0.3)$ ,  $\sigma = (1.8, 1.2, 0.8, 4.0)$

The behaviour of the estimations for all algorithms are the same that we explain in section 1.1.2. The major difference is that the Q-values estimations for all algorithms, when exploring takes more time, because the standard deviation is larger. Because of that it is harder for the learner to estimate the correct Q-values.



### 1.2.3 Histograms

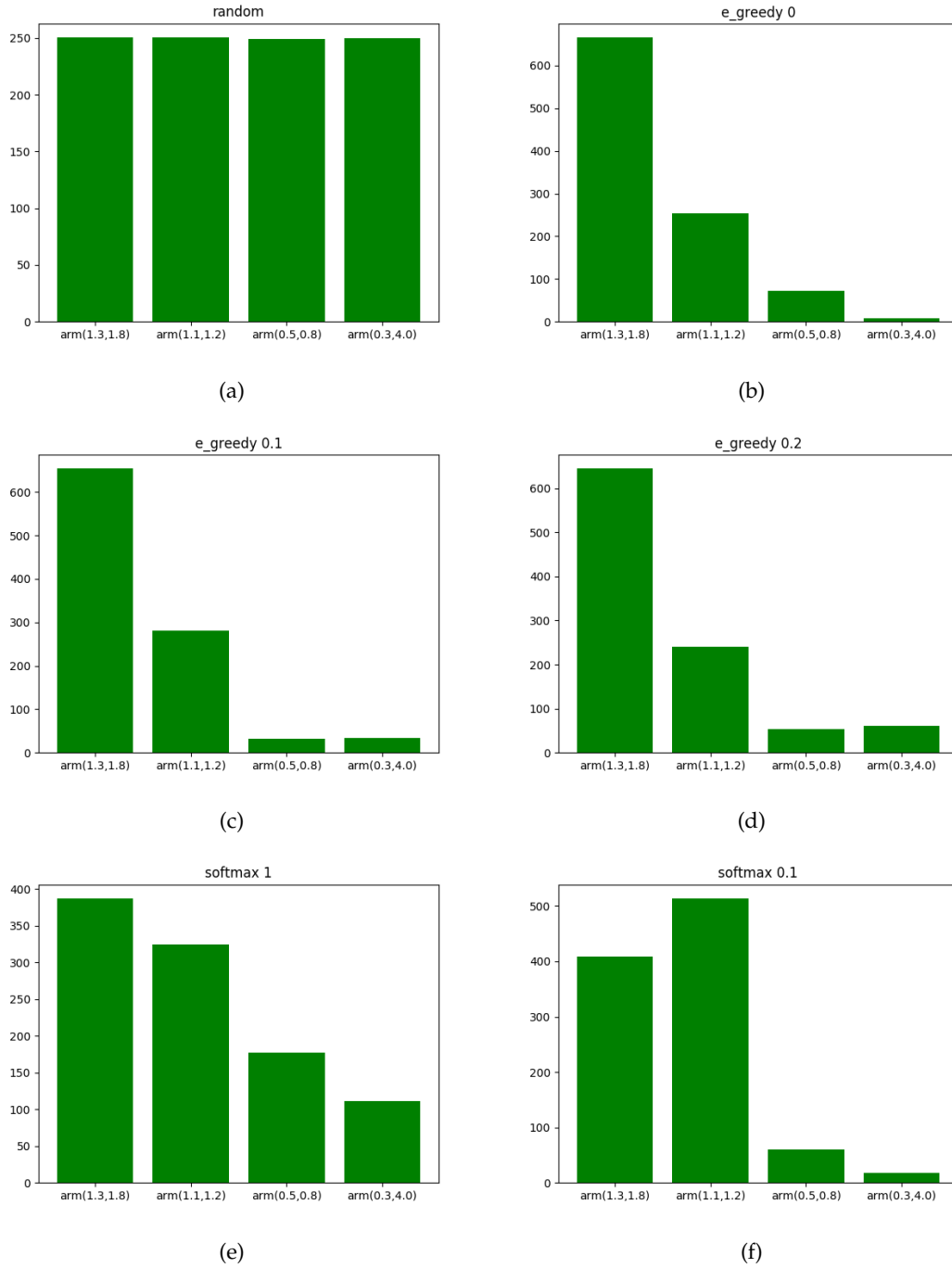


Figure 6: Histograms showing the number of times each action is selected per selection strategy with  $\mu = (1.3, 1.1, 0.5, 0.3)$ ,  $\sigma = (1.8, 1.2, 0.8, 4.0)$

The actions selection behaviour for all algorithms seems to follow the same pattern that we explain in section 1.1.3. The major difference is that for the one that exploits the most, for example 0-greedy, the selection of the *worst* action have a higher bars. This is because of the larger standard deviation for all actions. It is harder to determine which one is the optimal selection *i.e.* the exploration at the beginning is necessary, because the Q-values are estimated delayed compared to when the standard deviation is smaller.

### 1.3 Exercise 3

#### 1.3.1 Average rewards

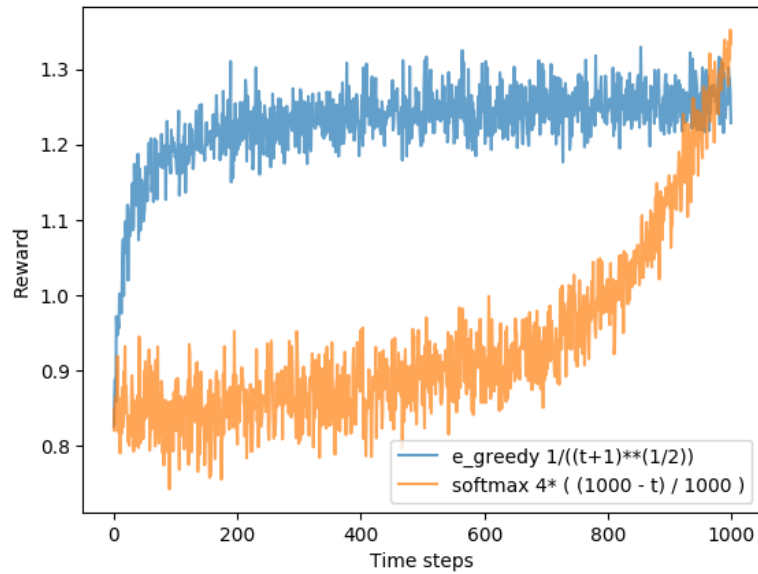


Figure 7: Average rewards for all algorithms

Dynamic  $\epsilon$ -greedy doesn't seem to perform better than the static ones. The resulting rewards tend to be identical. Observing the graph shows that it takes more time to stabilize to the best arm.

Dynamic softmax seems to be better, if you consider the long term reward. It takes to the 1000 epoch for the algorithm to find the best action reward in opposition to the softmax using a temperature of 0.1, where it is close to the best arm at around epoch 200. But not as close as the dynamic softmax, where it is closer to the best Q-value, when it is stabilized.

#### 1.3.2 Q-values estimation

Figure 8 shows the Q-values estimation overtime of the dynamic selection methods  $\epsilon$ -greedy and softmax.

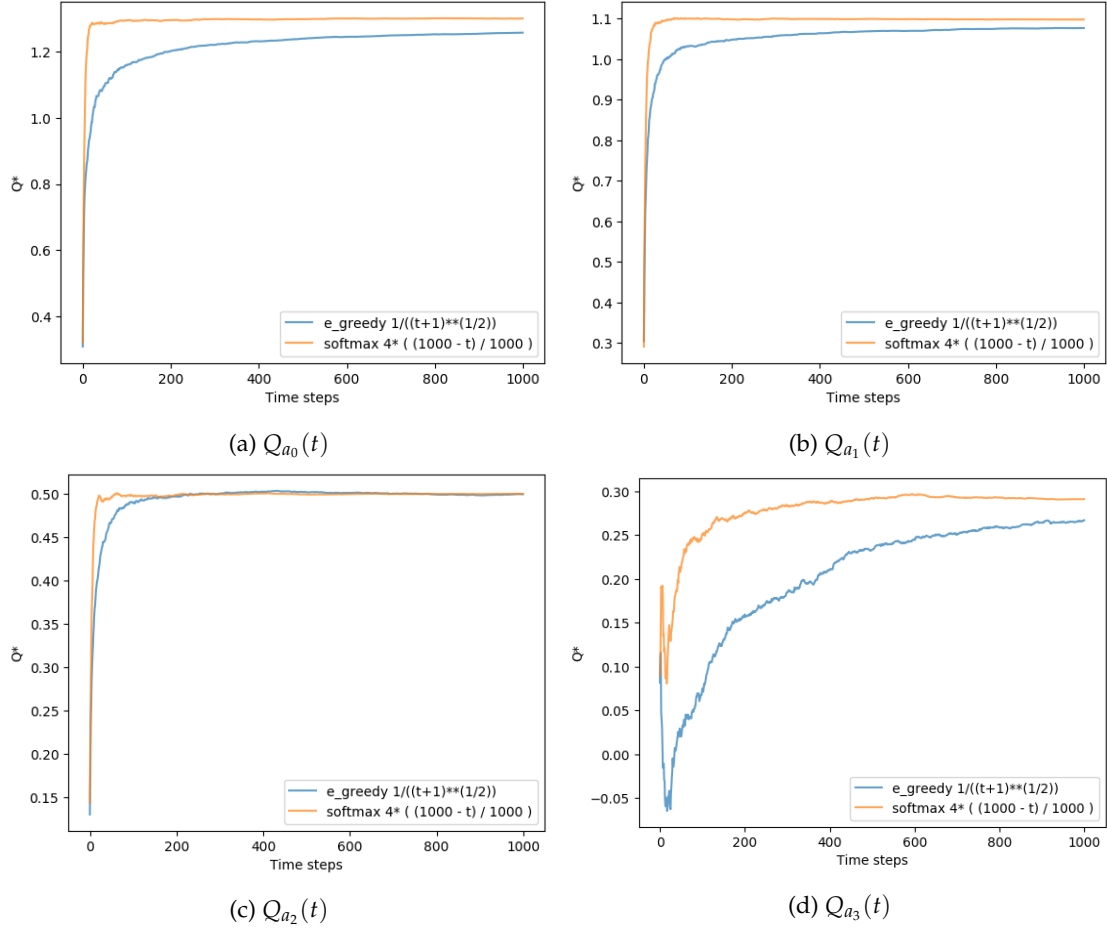


Figure 8: Plot per arm showing the  $Q_{a_i}^*$  of that action along with the actual  $Q_{a_i}$  estimate over time

**Dynamic  $\epsilon$ -greedy** Since the dynamic e-greedy method is randomized at the beginning and more greedy over time, the Q-values are correctly estimated and tend to the correct  $Q^*$  over time.

**Dynamic softmax** We have a large temperature at the beginning and a smaller temperature over time. This leads to an algorithm that explores a lot at the beginning, therefore estimating the Q-values very fast.

### 1.3.3 Histograms

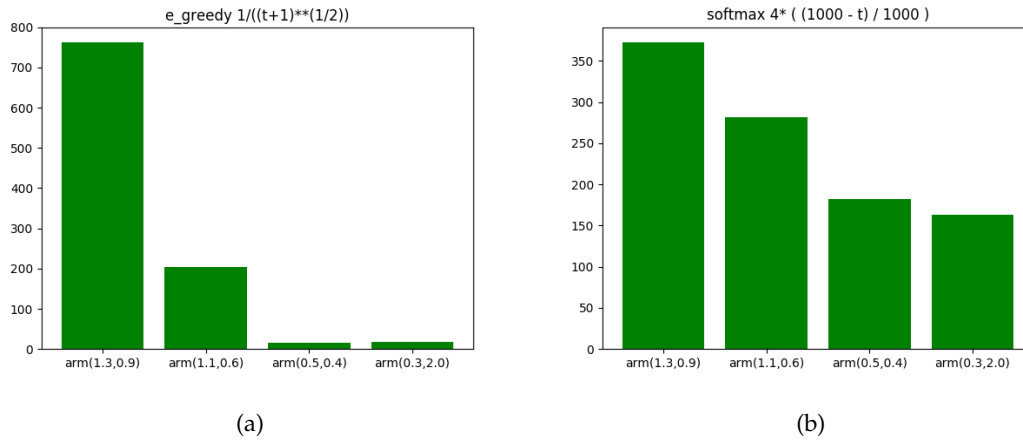


Figure 9: Histograms showing the number of times each action is selected per selection strategy

**Dynamic  $\epsilon$ -greedy** The action selection tends to the optimal action #0 very fast. This is because that the probability  $\epsilon$  drops very fast overtime, which lead to the greedy selection behaviour.

**Dynamic softmax** Since dynamic softmax has a large temperature when  $t$  is small, therefore explore a lot, we can see that all actions are selected. Overtime to selection will be more greedy, and since the Q-values are accurate because of the exploration, when  $t$  is high enough, the action #0 and #1 will be selected more often. That explains the *stairs* form of the histogram.

## 2 Climbing game

When using the formula, it is important to have a big tau, because  $\exp(Q/\tau)$  will give errors if Q is too big. (See how EV(a) is calculated)

I also remarked that when using EV with  $\max_{\tau} \text{ewardsormax}_q \text{when min}_t \tau = 0.001$  the 11 is good, but when 0.1 or 1, not.