# CS6700 : Reinforcement Learning

## Programming Assignment 1

Akshit Kumar

*EE14B127*

24th February 2017

# Contents

# 1   Introduction

## 1.1   Goal

The goal of this assignment is to evaluate differenet bandit algorithms on an N-armed bandit testbed. This report analyses and presents the results obtained by simulating an N-armed bandit problem on 3 types of Action-Value Methods:

- $\epsilon$-greedy

- Sampling from softmax distribution

- UCB-1 Algorithm

Each of the above algorithms are first run on a test-bed of 10-armed bandits with the results averaged over 2000 different bandit problems for 1000 time steps. In addition to this, we also compare the three algoritms on a test-bed of 1000 arms and see how the algorithm performs for the three algorithms perform by testing it on different time-steps.

## 1.2   Implementation Details

### 1.2.1   True Expected Payoff of Arms

For the purposes of experimental analysis, the true expected payoff of the arms for different bandit problems is sampled from a standard normal distribution offset by a true reward value. Therefore the true expected payoff for each arm is

$$q_*(a) = TrueReward + \mathcal{N}(0,1)$$

### 1.2.2   Initialization of Estimates of Expectation of Arms

Initially all the estimates of the expectation of arms is set to zero. Therefore we have

$$Q_0(a) = 0 \qquad \forall a \in \mathcal{A}$$

### 1.2.3   Sampling of Rewards for Arms

The rewards for each of the arms at each time step in each of the bandit problems is sampled from a standard normal distribution offset by the true expected payoff of that arm. Therefore we have

$$R_i = q_*(a) + \mathcal{N}(0,1)$$

where $R_i$ is the reward at time step $i$ and $\mathcal{N}(0,1)$ adds noise to the true expected payoff for that action.

### 1.2.4 Update of the Estimate of Arms

For updating the estimates of the arm, incremental updates are made and the update equation is given as

$$Q_{t+1} = Q_t + \frac{1}{t}[R_t - Q_t]$$

# 2 Action-Value Algorithms

In this section we discuss the algorithms of the 3 types of action-value methods mentioned above.

## 2.1 $\epsilon$-Greedy

---
**Algorithm 1** $\epsilon$-greedy Algorithm

---
1: **procedure** $\epsilon$-GREEDY PROCEDURE
2:     **for** each action $a \in \mathcal{A}$ **do**
3:         Q(a)$\leftarrow$0
4:         N(a) $\leftarrow$ 0
5:     **end for**
6:     **loop**
7:         **if** sample from $\mathcal{N}(0,1) > \epsilon$ **then**
8:             $\mathcal{A} \leftarrow \mathrm{argmax}_a Q(a)$
9:         **else**
10:             $\mathcal{A} \leftarrow \mathrm{randint}(0, \# \text{ of arms})$
11:         **end if**
12:         $\mathcal{R} \leftarrow getReward(\mathcal{A})$
13:         $N(\mathcal{A}) \leftarrow N(\mathcal{A}) + 1$
14:         $Q(\mathcal{A}) \leftarrow Q(\mathcal{A}) + \frac{1}{N(A)}[R - Q(A)]$
15:     **end loop**
16: **end procedure**

---

This algorithm introduces some randomness in selection of action to take at step $t+1$ which allows us to explore more actions and not settle for a suboptimal arm. At each step, with probability $\epsilon$ , a random exploratory move is played and with probability $1 - \epsilon$, a greedy action selection is made based on the current estimates of the expectation of the arms.

## 2.2 Sampling from softmax distribution

---
**Algorithm 2** Sampling from Softmax Distribution

---
    **procedure** SAMPLING FROM SOFTMAX DISTRIBUTION
2:     **for** each action $a \in \mathcal{A}$ **do**
        Q(a) $\leftarrow$ 0
4:         N(a) $\leftarrow$ 0
    **end for**
6:     **loop**
        $\mathcal{A} \sim \{\frac{e^{Q_t(a_i)/\tau}}{\sum_a e^{Q_t(a)/\tau}}\}$
8:         $\mathcal{R} \leftarrow getReward(\mathcal{A})$
        $N(\mathcal{A}) \leftarrow N(\mathcal{A}) + 1$
10:         $Q(A) \leftarrow Q(A) + \frac{1}{N(A)}[R - Q(A)]$
    **end loop**
12: **end procedure**

---

In this algorithm, we sample from a softmax distribution which is given by

$$P_t = \frac{e^{Q_t(a_i)/\tau}}{\sum_a e^{Q_t(a)/\tau}}$$

where $\tau$ is the temperature parameter.

## 2.3 Upper Confidence Bound

---
**Algorithm 3** Upper Confidence Bound

---
 1: **procedure** UCB-1
 2:    **for** each action $a \in \mathcal{A}$ **do**
 3:       play action $a$
 4:       $\mathcal{R} \leftarrow getReward(\mathcal{A})$
 5:       $Q(a) \leftarrow \mathcal{R}$
 6:    **end for**
 7:    **loop**
 8:       $\mathcal{A} \leftarrow \text{argmax}_a\, Q(a) + \sqrt{\frac{2lnN}{T_i(n)}}$
 9:       $\mathcal{R} \leftarrow getReward(\mathcal{A})$
10:       $N(\mathcal{A}) \leftarrow N(\mathcal{A}) + 1$
11:       $Q(A) \leftarrow Q(A) + \frac{1}{N(A)}[R - Q(A)]$
12:    **end loop**
13: **end procedure**

---

In this algorithm, we initially play all the actions once and then we take the subsequent actions in the remaining timesteps by being greedy with respect to the upper confidence bound on the current estimates of the actions. Therefore at each time step we choose,

$$a = \underset{a}{\text{argmax}}\, Q(a) + \sqrt{\frac{2lnN}{T_i(N)}}$$

# 3 Experimental Results

## 3.1 Plots for $\epsilon$-greedy Action Selection for 10 arms



Figure 1: Plot of average reward of $\epsilon$-greedy action-value methods on the 10-armed testbed. These data are averages over 2000 runs with different bandit problems.



Figure 2: Plot of % of optimal action selection of $\epsilon$-greedy action-value methods on the 10-armed testbed. These data are averages over 2000 runs with different bandit problems.

## 3.2 Plots for Softmax Action Selection for 10 arms



Figure 3: Plot of average reward of softmax action selection methods on the 10-armed testbed. These data are averages over 2000 runs with different bandit problems.
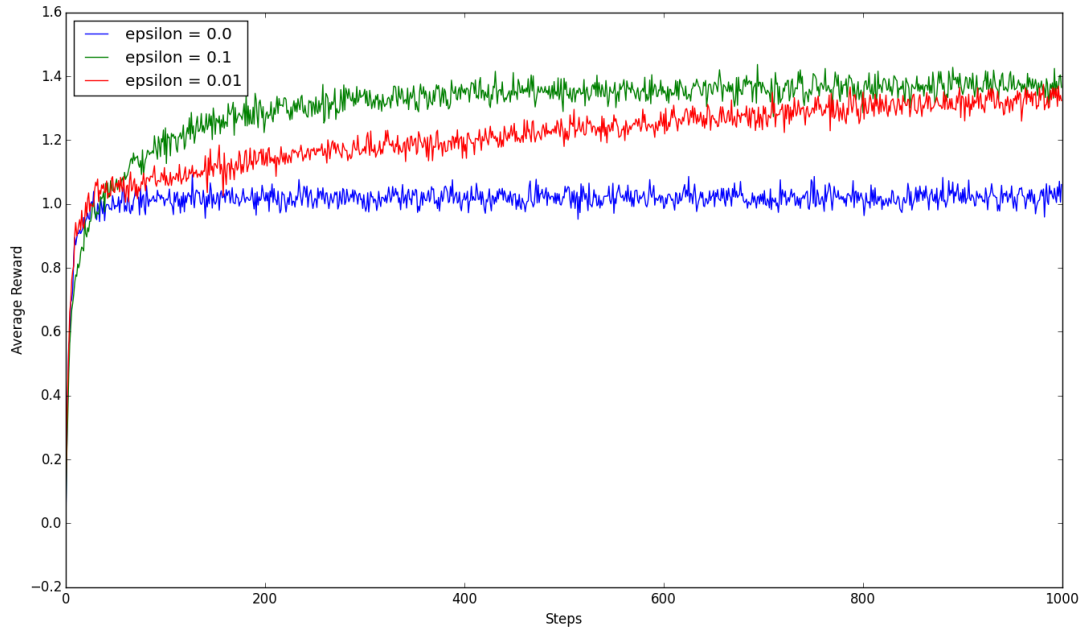


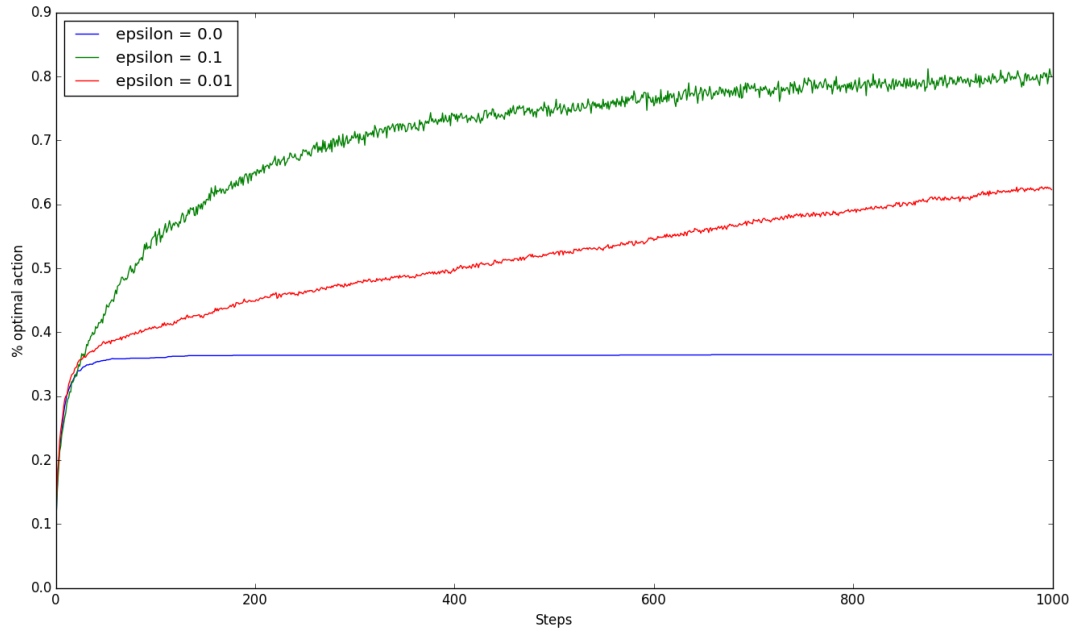Figure 4: Plot of average reward of softmax action selection methods on the 10-armed testbed. These data are averages over 2000 runs with different bandit problems.

## 3.3 Plots for UCB vs $\epsilon$-greedy vs Softmax Action Selection for 10 arms



Figure 5: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 10-armed testbed over 1000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.



Figure 6: Plot comparing the % optimal action selection using UCB, $\epsilon$-greedy and Softmax Action Selection on the 10-armed testbed over 1000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.

## 3.4 Plots for UCB vs $\epsilon$-greedy vs Softmax Action Selection for 1000-arms

### 3.4.1 Plots for 1000 time-steps



Figure 7: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 1000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.

Figure 8: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 1000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.
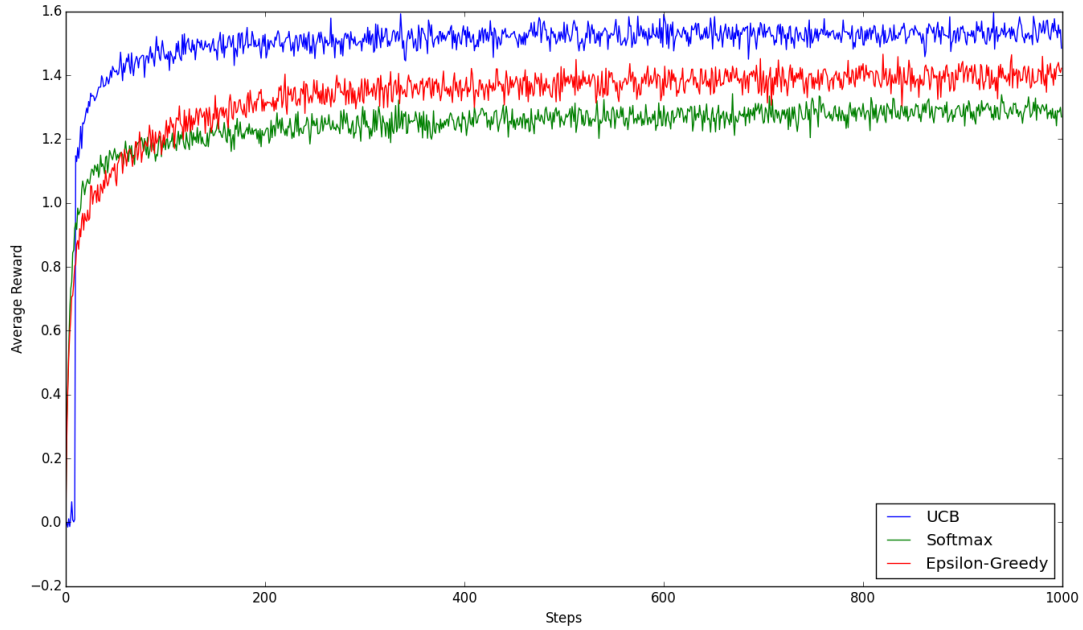
### 3.4.2 Plots for 10000 time-steps



Figure 9: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 10000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.



Figure 10: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 10000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.
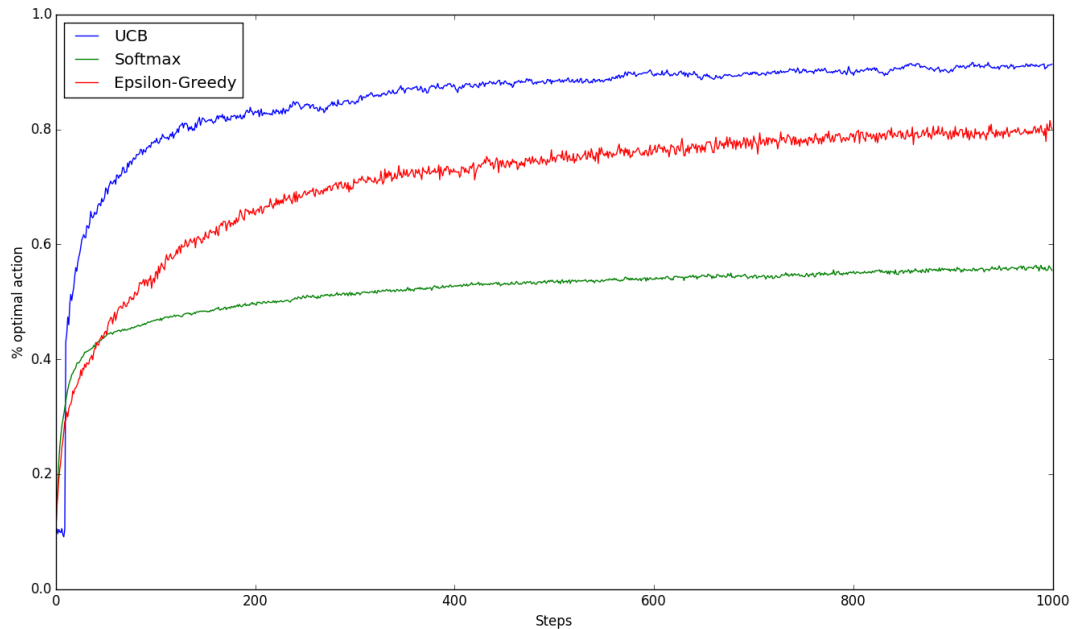
### 3.4.3 Plots for 20000 time-steps



Figure 11: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 20000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.
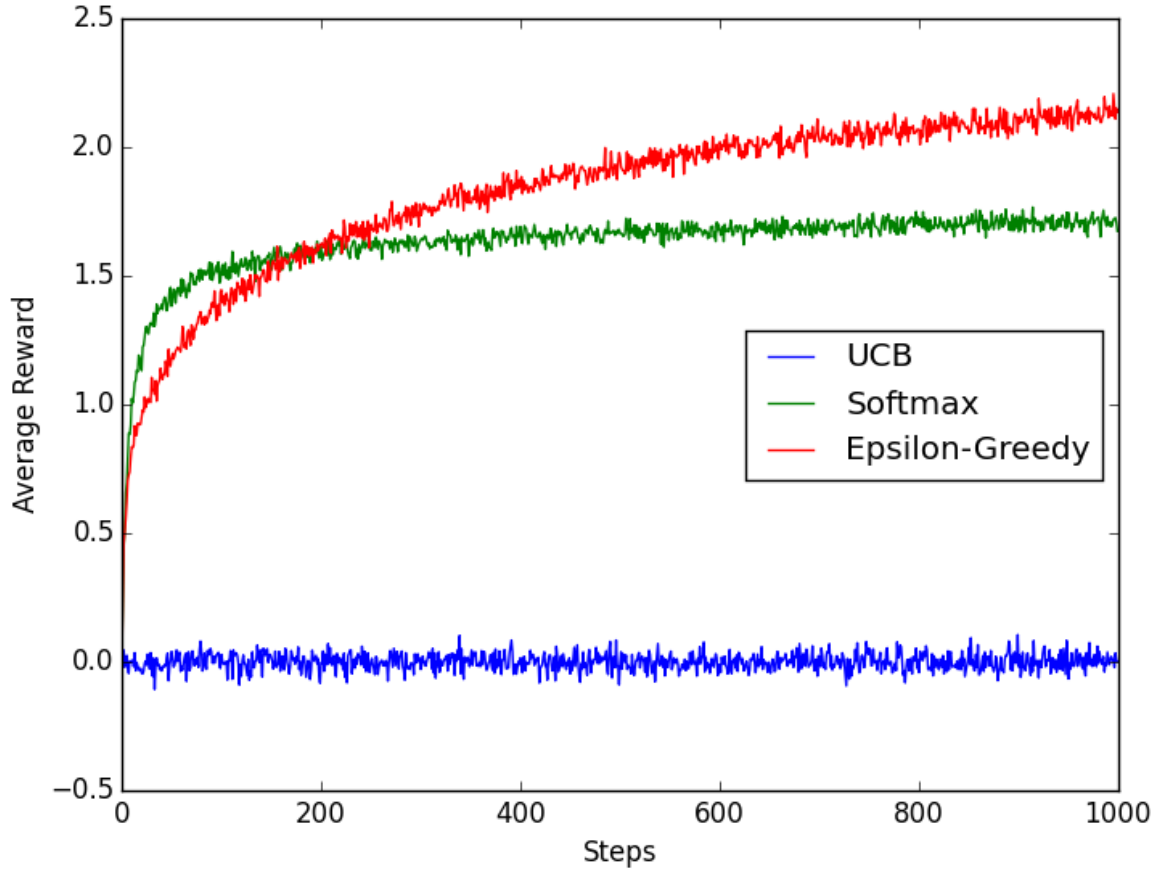


Figure 12: Plot comparing the average reward accrued using UCB, $\epsilon$-greedy and Softmax Action Selection on the 1000-armed testbed over 20000 timesteps. These data are averaged over 2000 runs of different bandit problems. For $\epsilon$-greedy approach, $\epsilon = 0.1$ and for softmax action selection, $\tau = 0.1$ where $\tau$ is the temperature parameter.
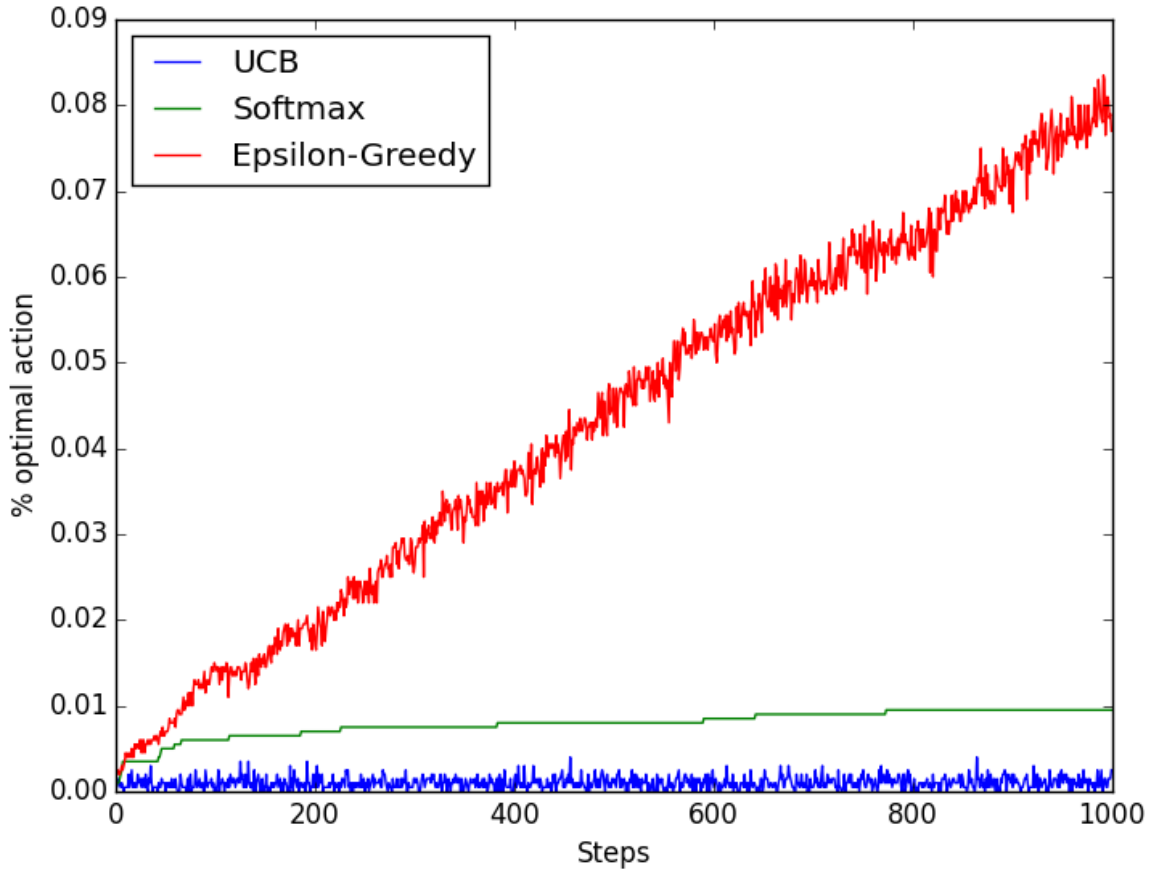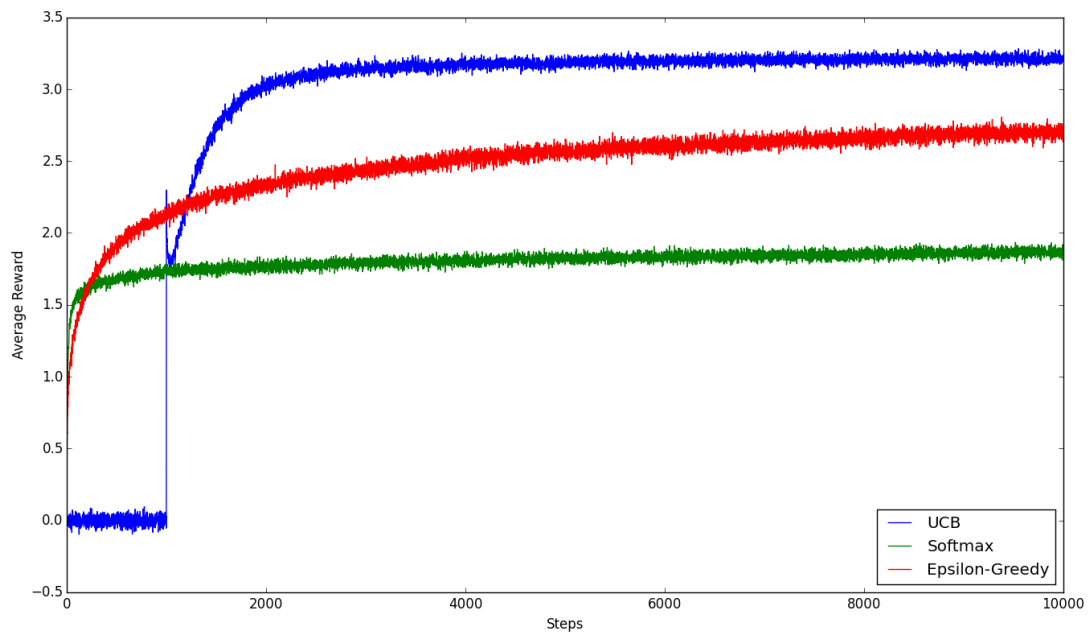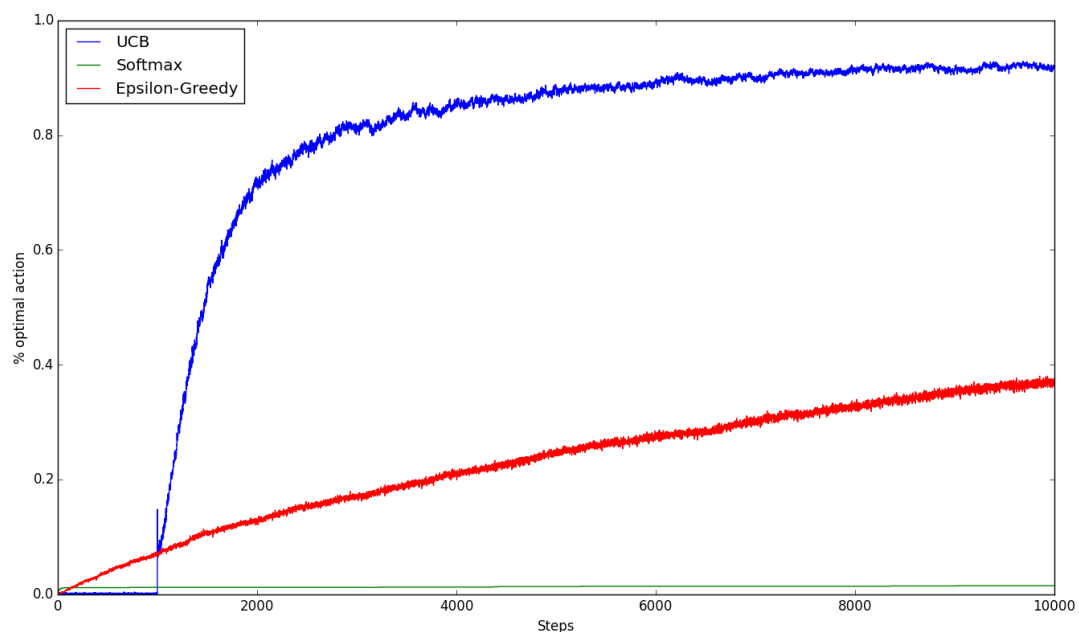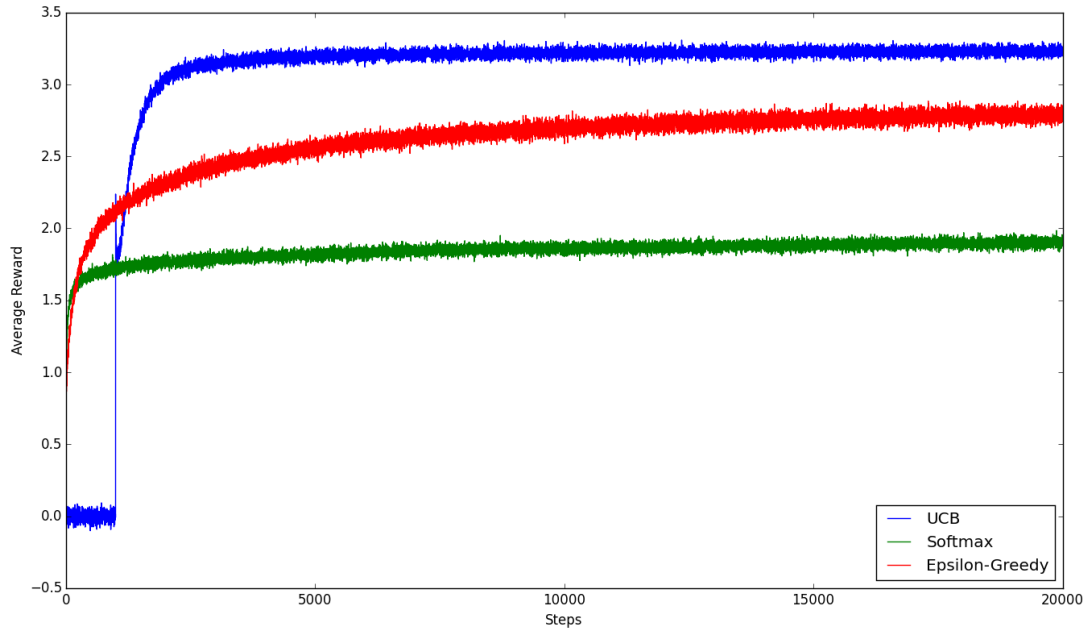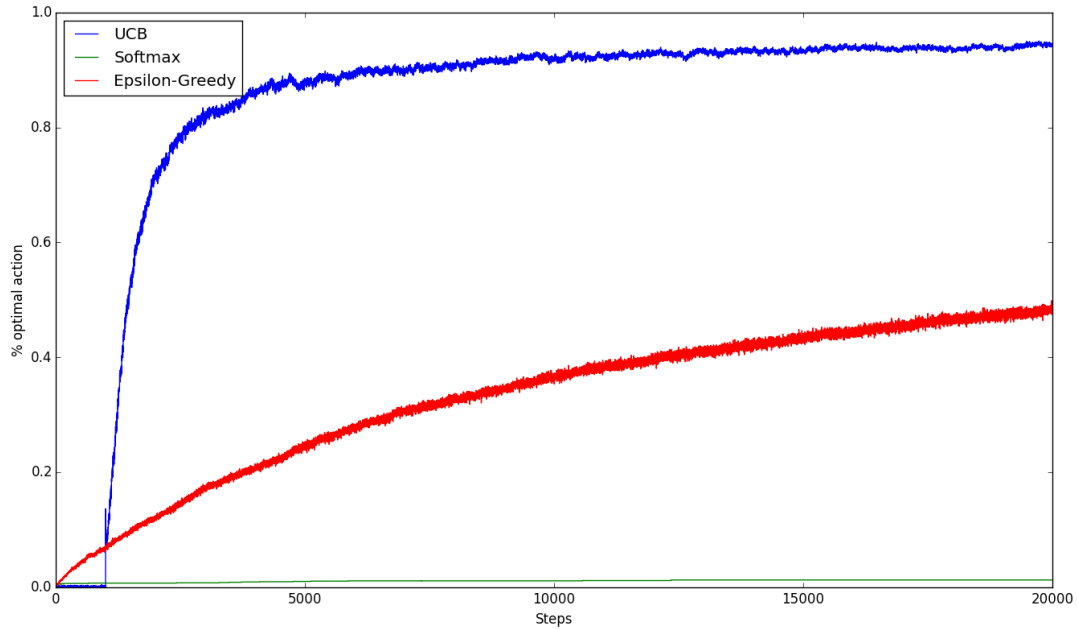
# 4 Analysis of the Plots

## 4.1 Plots for $\epsilon$-greedy action selection for 10 arms

The simulation for $\epsilon$-greedy action selection is run on a 10-armed bandit testbed for three different values of $\epsilon$ which are 0.0 (greedy) , 0.1 (being greedy 90 % of the time) , 0.01 (being greedy 99 % of the time). We can notice in the plot of these simulations that over a time-step of 1000, $\epsilon = 0.1$ does the best both in terms of average reward and % optimal action selection. Notice that in the long run $\epsilon = 0.01$ might do better than $\epsilon = 0.1$ both in terms of average reward and % optimal action selection. The greedy method ($\epsilon = 0.0$) performs significantly worse in the long run because it often gets stuck performing a suboptimal action. The % optimal action selection shows that the greedy action selection found the optimal action in only approximately one-third of the tasks. Therefore in two-third of the tasks,its initial samples of optimal action were bad and it never returned to them. In comparison to greedy, $\epsilon$-greedy perform better because they continue to explore and improve their chance of recognizing the optimal action.

## 4.2 Plots for Softmax Action Selection for 10 arms

The simulation for softmax action selection is run on a 10-armed bandit testbed for four different values of temperature which are 0.01, 0.1, 1 and 10. We can notice in the plot of these simulations that over a time-step of 1000, temperature of 0.1 does the best both in terms of average reward and % optimal action selection. We can notice that if the value of temperature is too low (ie 0.01), we essentially end up behaving greedily but it often gets stuck performing a suboptimal action. Similarly we can also notice that if the value of temperature is too high (ie 10), we essentially end up being random in our action selection and perform badly and hence % optimal action is only 10% ie we are picking up the optimal arm uniformly randomly. We can notice that temperature value of 0.1 is the most apt as it carefully balances exploration and exploitation. Also this is the only temperature which continues to grow for both the plots while the other temperatures seem to saturate out.

## 4.3 Plots for UCB vs $\epsilon$-greedy vs Softmax Action Selection for 10 arms

The simulation for comparison of UCB with other action selection plots is run on a 10-armed bandit testbed for $\epsilon = 0.1$ for $\epsilon$-greedy and temperature $= 0.1$ for softmax action selection over a 1000 timesteps. We can clearly notice that UCB does the best both in terms of average reward and % optimal action selection. Initially, UCB performs bad but that can be attributed to the fact that initially UCB plays all the onces once and is forced to take suboptimal actions. Also in comparison $\epsilon$-greedy performs better that softmax action selection for this particular choice of $\epsilon$ and temperature. But we can conclude that UCB performs the best in case of 10 arms over 1000 timesteps.

## 4.4 Plots for UCB vs $\epsilon$-greedy vs Softmax Action Selection for 1000 arms

Now we do the comparison of algorithms as the number of arms increase. We increase the number of arms from 10 to 1000 and compare the algorithms for different values of time-steps.

### 4.4.1 Plots for 1000 timesteps

If we keep the number of timesteps the same as the one in the previous question, we can notice that UCB performs the worst in comparison to $\epsilon$-greedy and softmax action selection. $\epsilon$-greedy performs the best eventually(ie in 1000 timesteps) because it starts picking up sub-optimal action which still does better than UCB because UCB is continuing to play each arm one and hence doesn't repeatedly take a sub-optimal action. In the $\epsilon$-greedy action selection doesn't visit all the arms and selects action from a set of few sampled arms. It is not fair to compare $\epsilon$-greedy and softmax action selection amongst themselves as these are parameterized by $\epsilon$ and temperature

and hence for different values of $\epsilon$ and temperature may yield different results but this is sufficient to conclude that both $\epsilon$-greedy and softmax action selection perform better than UCB for 1000 time-steps simply because UCB continues to play each arm once while $\epsilon$-greedy and softmax action selection behave greedily majority of the time. This can also be seen from the %optimal action plot which shows that UCB selects optimal action once in 1000 times which should be expected. This is also not fair on UCB and hence in the next analysis we increase the time step from 1000 to 10000 and 20000.

### 4.4.2 Plots for 10000 timesteps

As we increase the number of timesteps from 1000 to 10000 we can immediately see that UCB performs best in comparison to $\epsilon$-greedy and softmax action selection. Initially again, UCB performs the worst but that can again be attributed to the fact that UCB is playing each arm once. We can also see that there is sudden shoot in average reward after 1000 steps which can be attributed to the fact that UCB now starts being greedy with respect to the upper confidence bound on the estimates. Again we see that $\epsilon$-greedy does better than softmax action selection which is not a fair comparison as they are parameterised by $\epsilon$ and temperature.

### 4.4.3 Plots for 20000 timesteps

The results for 20000 timesteps are similar to the ones for 10000 timesteps. One thing to notice though is that $\epsilon$-greedy seems to be converging to UCB in terms of average reward and % optimal action selection.

## 5  Conclusion

- $\epsilon$-greedy T = 1000 and N = 10

  - maximum average reward is obtained for $\epsilon = 0.1$
  - maximum optimal action is obtained for $\epsilon = 0.1$

- softmax action selection T = 1000 and N = 10

  - maximum average reward is obtained for temperature = 0.1
  - maximum optimal action is obtained for temperature = 0.1

- UCB vs $\epsilon$-greedy vs softmax action selection T = 1000 and N = 10

  - maximum average reward is obtained for UCB
  - maximum optimal action is obtained for UCB

- UCB vs $\epsilon$-greedy vs softmax action selection T = 1000 and N = 1000

  - maximum average reward is obtained for $\epsilon$-greedy with $\epsilon = 0.1$
  - maximum optimal action is obtained for $\epsilon$-greedy with $\epsilon = 0.1$

- UCB vs $\epsilon$-greedy vs softmax action selection T = 10000 and N = 1000

  - maximum average reward is obtained for UCB
  - maximum optimal action is obtained for UCB

- UCB vs $\epsilon$-greedy vs softmax action selection T = 20000 and N = 1000

  - maximum average reward is obtained for UCB
  - maximum optimal action is obtained for UCB

# 6 Code Listing

## 6.1 Bandit Class for object instantiation

```python
1  # Import necessary libraries
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  ''' Define a Bandit class which has the following parameters
6  @param numArm : number of actions in the K-armed bandit problem
7  @param epsilon : Set the epsilon value for epsilon-greedy action selection. Default
       value is 0.0
8  @param temperature : Set the temperature parameter of softmax action selection.
       Default value is 1.0
9  @param epsilonGreedy : Flag to select that action-selection is based on epsilon-greedy
       algorithm. Default value is False
10 @param softmax : Flag to select that action-selection is based on the softmax
       algorithm. Defalut value is False
11 @param ucb : Flag to select that action-selection is based on the softmax algorithm.
       Default value is False
12 @param initial_estimate : Setting the initial estimates for estimates all the arms
13 @param trueReward : Setting the true reward for the all arms
14 '''
15
16 class Bandit:
17     # Constructor function to initialise the parameters for the bandit instantiation
18     def __init__(self, numArm = 10, epsilon = 0.0, temperature = 1.0, epsilonGreedy =
           False, softmax = False, ucb =
19     False, initial_estimate = 0.0, trueReward = 0.0):
20         # Initialisation of parameters
21         self.numArm = numArm
22         self.epsilon = epsilon
23         self.temperature = temperature
24         self.time = 0
25         self.averageReward = 0
26         self.epsilonGreedy = epsilonGreedy
27         self.softmax = softmax
28         self.ucb = ucb
29         self.armIndices = np.arange(self.numArm)
30         self.qTrue = trueReward * np.ones(self.numArm) # True rewards for each action
31         self.qEst = initial_estimate *  np.ones(self.numArm) # Estimate of the reward
               for each action
32         self.actionCount = np.zeros(self.numArm) # Action Selection count
33         # Get the true rewards for each action
34         for i in range(0,self.numArm):
35             self.qTrue[i] = self.qTrue[i] + np.random.randn()
36         self.bestAction = np.argmax(self.qTrue)
37
38     # Function to choose the action based on what algorithm is being followed
39     def chooseAction(self):
40         if self.epsilonGreedy == True:
41             armToPlay = np.random.randint(0,self.numArm) # Choose a random action
42             # With probability greater than epsilon, choose the action greedily
43             if np.random.random() > self.epsilon:
44                 armToPlay = np.argmax(self.qEst)
45             return armToPlay
46
47         if self.softmax == True:
48             expo_val = np.exp(self.qEst / self.temperature)
49             prob_val = expo_val / np.sum(expo_val) # Get the probability distribution
                   to be sampled from
50             rnd = np.random.random() # Sample a random number
51             # Use the CDF method to sample from the Gibbs distribution
52             for i,w in enumerate(prob_val):
53                 rnd -= w
54                 if rnd < 0:
```

```
55                         return i
56                 return 0
57
58        if self.ucb == True:
59            # play each action once at first
60            if(np.count_nonzero(self.actionCount) < self.numArm):
61                return np.argmin(self.actionCount)
62            # after having played each arm once, choose action greedily using the
                upper confidence bound
63            return np.argmax(self.qEst + np.sqrt(np.log(self.time + 1) /
                (np.asarray(self.actionCount) + 1)))
64
65    # Function to get the reward on playing the action, increment the count of action
          and update the estimate of that action
66    def getReward(self,action):
67        reward = np.random.randn() + self.qTrue[action]
68        self.time += 1
69        self.averageReward = (self.time - 1.0) / self.time * self.averageReward +
              reward / self.time # keep a running average of the reward
70        self.actionCount[action] += 1 # increase the action count
71        self.qEst[action] += 1.0 / self.actionCount[action] * (reward -
              self.qEst[action]) # update the estimate of the action
72        return reward
```
../bandit.py

## 6.2 Python Library for Simulating the Bandit Problem

```
1  # Import the necessary libraries
2  import numpy as np
3
4  ''' Function to simulate the bandit problem
5  @param numDiffBandits : variable stating the number of different bandit problems
6  @param timeSteps : variable stating the number of time steps over which the actions
       are taken
7  @param bandits : list contains the different bandit problems
8  '''
9  def simulateBandit(numDiffBandits, timeSteps, bandits):
10     optimalActionCount = np.zeros(shape=(len(bandits),timeSteps)) # 2D matrices for
           storing the count of optimal action selection in each time step for all the
           different bandit problems
11     averageReward = np.zeros(shape=(len(bandits),timeSteps)) # 2D matrices for storing
           the average reward in each time step for all the different bandit problems
12     for banditIndex, bandit in enumerate(bandits): # iterating through the all the
           different bandit problems
13         for i in range(0,numDiffBandits): # iterating through the number of different
               bandits
14             for j in range(0,timeSteps): # iterating through the time steps for each
                   bandit problem
15                 action = bandit[i].chooseAction() # select an action to take at each
                       time step depending on the action selection method for the
                       particular bandit
16                 reward = bandit[i].getReward(action) # get a reward for selecting that
                       action for that particular time step
17                 averageReward[banditIndex][j] += reward # calculate the cummulative
                       reward
18                 if action == bandit[i].bestAction:
19                     optimalActionCount[banditIndex][j] += 1 # increment the count of
                           action by 1
20         optimalActionCount[banditIndex] /= numDiffBandits # average out the optimal
               action count
21         averageReward[banditIndex] /= numDiffBandits # average out the cummulative
               reward
22     return optimalActionCount, averageReward
```
../simulateBandit.py

## 6.3 Python Script to Generate Plots for $\epsilon$-greedy action selection

```python
1  ''' Run this python script to generate the plots for the epsilon-greedy plot on a
       N-armed bandit testbed'''
2  from bandit import Bandit # import the bandit class to instantiate bandit objects for
       different bandit problems
3  from simulateBandit import simulateBandit # import the simulateBandit function to do
       the simulation of arm pulls and getting rewards
4  import matplotlib.pyplot as plt # for plotting graphs
5
6  ''' Function used for plotting the average reward and optimal action selection using
       epsilon greedy action selection methods
7  @param numDiffBandits : number of different bandit problems
8  @param timeSteps : number of different time steps over which the algorithm is run
9  '''
10 def epsilonGreedy(numDiffBandits, timeSteps):
11     epsilons = [0.0 , 0.1, 0.01] # taking three different values of epsilon
12     bandits = list() # making a list of bandit problems
13     # iterating through the list of epsilons and adding different bandit problems
14     for epsilonIndex, epsilon in enumerate(epsilons):
15         bandits.append([Bandit(epsilon = epsilon, epsilonGreedy = True) for _ in
               range(0,numDiffBandits)])
16     optimalActionCount, averageReward = simulateBandit(numDiffBandits, timeSteps,
           bandits)
17     # plotting the results for optimal action
18     plt.figure(0)
19     for epsilon, counts in zip(epsilons, optimalActionCount):
20         plt.plot(counts, label='epsilon = ' + str(epsilon))
21     plt.xlabel('Steps')
22     plt.ylabel('% optimal action')
23     plt.legend(loc='best')
24     # plotting the results for average reward
25     plt.figure(1)
26     for epsilon, rewards in zip(epsilons, averageReward):
27         plt.plot(rewards, label='epsilon = ' + str(epsilon))
28     plt.xlabel('Steps')
29     plt.ylabel('Average Reward')
30     plt.legend(loc='best')
31
32 epsilonGreedy(2000,1000) # running the simulation on 10-armed bandit testbed for 2000
       different bandit problems over 1000 time steps
33 plt.show()
```

../epsilonGreedy.py

## 6.4 Python Script to Generate Plots for softmax action selection

```python
1  ''' Run this python script to generate the plots for the softmax action selection
       plots on a N-armed bandit testbed'''
2  from bandit import Bandit # import the bandit class to instantiate bandit objects for
       different bandit problems
3  from simulateBandit import simulateBandit # import the simulateBandit function to do
       the simulation of arm pulls and getting rewards
4  import matplotlib.pyplot as plt # for plotting graphs
5
6  ''' Function used for plotting the average reward and optimal action selection using
       softmax action selection methods
7  @param numDiffBandits : number of different bandit problems
8  @param timeSteps : number of different time steps over which the algorithm is run
9  '''
10 def softmax(numDiffBandits,timeSteps):
11     temperatures = [0.01,0.1,1,10]  # taking 4 different values of temperate for the
           softmax distribution
12     bandits = list() # making a list of bandit problems
13     # iterating through the list of temperatures and adding different bandit problems
14     for tempIndex, temp in enumerate(temperatures):
```

```
15          bandits.append([Bandit(temperature = temp, softmax = True) for _ in
                range(0,numDiffBandits)])
16      optimalActionCount, averageReward = simulateBandit(numDiffBandits, timeSteps,
            bandits)
17      # plotting the results for optimal action
18      plt.figure(0)
19      for temp, counts in zip(temperatures, optimalActionCount):
20          plt.plot(counts, label='temp = ' + str(temp))
21      plt.xlabel('Steps')
22      plt.ylabel('% optimal action')
23      plt.legend(loc='best')
24      # plotting the results for average reward
25      plt.figure(1)
26      for temp, rewards in zip(temperatures, averageReward):
27          plt.plot(rewards, label='temp = ' + str(temp))
28      plt.xlabel('Steps')
29      plt.ylabel('Average Reward')
30      plt.legend(loc='best')
31
32  softmax(2000,1000) # running the simulation on 10-armed bandit testbed for 2000
        different bandit problems over 1000 time steps
33  plt.show()
```

../softmax.py

## 6.5 Python Script to Generate Plots for comparing UCB vs $\epsilon$-greedy action selection vs softmax action selection

```
1  ''' Run this python script to generate the plots for comparing the UCB algorithm vs
       epsilon greedy and softmax on a N-armed bandit testbed'''
2  from bandit import Bandit # import the bandit class to instantiate bandit objects for
       different bandit problems
3  from simulateBandit import simulateBandit # import the simulateBandit function to do
       the simulation of arm pulls and getting rewards
4  import matplotlib.pyplot as plt # for plotting graphs
5
6  ''' Function used for plotting the average reward and optimal action selection using
       ucb,epsilon greedy and softmax action selection methods
7  @param numArms : number of arms in the testbed
8  @param numDiffBandits : number of different bandit problems
9  @param timeSteps : number of different time steps over which the algorithm is run
10 '''
11 def ucb(numArms,numDiffBandits, timeSteps):
12     algos = ['UCB','Softmax','Epsilon Greedy'] # iterating through the different
           algorithms
13     bandits = [[],[],[]]
14     bandits[0] = [Bandit(numArm = numArms, ucb = True) for _ in
           range(0,numDiffBandits)] # list of bandits following ucb
15     bandits[1] = [Bandit(numArm = numArms, softmax = True, temperature = 0.1) for _ in
           range(0,numDiffBandits)] # list of bandits following softmax action selection
16     bandits[2] = [Bandit(numArm = numArms, epsilonGreedy = True, epsilon = 0.1) for _
           in range(0,numDiffBandits)] # list of bandits following epsilon greedy selection
17     optimalActionCount, averageReward = simulateBandit(numDiffBandits, timeSteps,
           bandits)
18     # plotting the results for optimal action
19     plt.figure(0)
20     plt.plot(optimalActionCount[0],label='UCB')
21     plt.plot(optimalActionCount[1],label='Softmax')
22     plt.plot(optimalActionCount[2],label='Epsilon-Greedy')
23     plt.xlabel('Steps')
24     plt.ylabel('% optimal action')
25     plt.legend(loc = 'best')
26     # plotting the results for average reward
27     plt.figure(1)
28     plt.plot(averageReward[0],label = 'UCB')
29     plt.plot(averageReward[1],label = 'Softmax')
```

```
30        plt.plot(averageReward[2],label = 'Epsilon-Greedy')
31        plt.xlabel('Steps')
32        plt.ylabel('Average Reward')
33        plt.legend(loc = 'best')
34
35  ucb(1000,2000,20000) # running the simulation on 1000-armed bandit testbed for 2000
        different bandit problems over 20000 time steps
36  plt.show()
```
<center>../ucb.py</center>