2024 - 1 Reinforcement Learning - Home work 1

1. Gittin Index: Measure of the reward that can be achieved by a sequence of actions from the present state onwords with the probability that it will be terminated in the future.

Gittin Index theorem helps on deciding a strategy of explore exploit trade-off by measuring the reward of certain action to give the best reward or discounted reward.

Thompson Sampling: Thompson sampling is an algorithm for online decision problems where actions are taken sequentially in a manner that must balance between exploiting what is known to maximize immediate performance and investing to accumulate new information that may improve future performance. It's particularly useful in situations where the decision has the outcomes of uncertainty.

Thompson sampling tends to explore actions with uncertainty outcomes while exploiting actions that are likely to yield high rewards, by balancing the trade off of exploration and exploitation.

2. In the first lo steps, the agents has explored the available actions which are the lo actions. The effect of selecting lo actions is that the action will be guaranteed to be picked at least once in the first lo steps.

Given the parameter $c > 0$, the agent is likely to prioritize the action with the highest return from the first sample. As of result, it may spike the reward/performance. As the agent try to do the other action, the agent will be much more forces / to do the explore option making it a sudden drop and try to get more reward as time moves on.

3.

$$Pr\{A_t = a\} = \frac{e^{H_t(a)}}{\sum_{d=1}^{h} e^{H_t(b)}} \qquad \text{sigmoid function} = \frac{1}{1+e^{-x}}$$

$$= \frac{e^{H_t(a)}}{e^{H_t(a)} + e^{H_t(b)}} = \frac{1}{1 + \frac{e^{H_t(b)}}{e^{H_t(a)}}}$$

$$= \frac{1}{1 + e^{H_t(b) - H_t(a)}}$$

if $x = H_t(b) - H_t(a)$.

then $= \dfrac{1}{1 + e^{-x}}$ ✳

4. Prove that :

$$H_{t+1}(A_t) = H_t(A_t) + a(R_t - \bar{R}_t)(1 - \pi_t(A_t)) \quad \&$$

$$H_{t+1}(a) = H_t(a) - a(R_t - \bar{R}_t)\pi_t(a) \qquad a \neq A_t$$

$$H_{t+1}(a) = H_t(a) + a\frac{\partial E[R_t]}{\partial H_t(a)} \qquad E[R_t] = \sum_x \pi_t(x)\, q_*(x)$$

$$= H_t(a) + a\left(\sum_x \left\{\frac{\partial}{\partial H_t(a)}[\pi_t(x) \cdot q_*(x)]\right\}\right)$$

$$= H_t(a) + a\left(\sum_x \left(q_*(x) \frac{\partial \pi_t(x)}{\partial H_t(a)}\right)\right)$$

$B_t$ : Baseline (we add Baseline to compare the current and previous value).

$$= H_t(a) + a\left(\sum_x (q_*(x) - B_t)\frac{\partial \pi_t(x)}{\partial H_t(a)}\right)$$

$$= H_t(a) + a\left(\underline{\sum_x (\pi_t(x)\cdot (q_*(x) - B_t)\frac{\partial \pi_t(x)}{\partial H_t(a)})}\,\Big/\,\pi_t(x)\right)$$

$$= H_t(a) + a\, E\left[(q_*(A_t) - B_t)\frac{\partial \pi_t(A_t)}{\partial H_t(a)}\,\Big/\,\pi_t(A_t)\right]$$

$$= H_t(a) + a\, E\left[(R_t - \bar{R}_t)\frac{\partial \pi_t(A_t)}{\partial H_t(a)}\,\Big/\,\pi_t(A_t)\right].$$

$$\frac{\partial \pi_t(x)}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \pi_t(x).$$

$$= \frac{\partial}{\partial H_t(a)} \left[ \frac{e^{H_t(x)}}{\sum_{y=1}^{u} e^{H_t(y)}} \right]$$

$$f'(x) = \frac{q'(x) h(x) - q(x) h'(x)}{h(x)^2}$$

$$= \frac{\frac{\partial e^{H_t(x)}}{\partial H_t(y)} \sum_{y=1}^{u} e^{H_t(y)} - e^{H_t(x)} \cdot e^{H_t(a)}}{\left( \sum_{y=1}^{u} e^{H_t(y)} \right)^2}$$

$$=) \quad \frac{\partial e^{H_t(x)}}{\partial H_t(y)} = \mathbb{1}_{a=x} \cdot e^{H_t(x)}$$

$$= \frac{\mathbb{1}_{a=x} e^{H_t(x)} \cdot \sum_{y=1}^{u} e^{H_t(y)} - e^{H_t(x)} \cdot e^{H_t(a)}}{\left( \sum_{y=1}^{u} e^{H_t(y)} \right)^2}$$

$$= \frac{\mathbb{1}_{a=x} e^{H_t(x)} \cdot \sum_{y=1}^{u} e^{H_t(y)}}{\left( \sum_{y=1}^{u} e^{H_t(y)} \right)^2} - \frac{e^{H_t(x)} e^{H_t(a)}}{\left( \sum_{y=1}^{u} e^{H_t(y)} \right)^2}$$

$$= \mathbb{1}_{a=x} \pi_t(x) - \pi_t(x) \pi_t(a)$$

$$= \pi_t(x) \left( \mathbb{1}_{a=x} - \pi_t(a) \right).$$

---

$$H_{t+1}(a) = H_t(a) + a \left[ \frac{(R_t - \bar{R}_t)(\pi_t(x)(\mathbb{1}_{a=x} - \pi_t(a))}{(\pi_t(x))} \right]$$

$$= H_t(a) + a \left[ (R_t - \bar{R}_t)(\mathbb{1}_{a=x} - \pi_t(a)) \right]$$

$$\text{if } \mathbb{1}_{a=x} = 1, \text{ then } a = x. \quad \Big| \text{ else} = 0.$$

$$H_{t+1}(A_t) = H_t(A_t) + a(R_t - \bar{R}_t)(1 - \pi_t(A_t)) \quad \Big| \quad H_{t+1}(a) = H_t(a) + a(R_t - \bar{R}_t)(-\pi_t(a))$$
$$= H_t(a) \cdot - a(R_t - \bar{R}_t)(\pi_t(a))$$

## 5a- add $l_1$

$$H_{t+1}(a) = H_t(a) + a \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial \left( \sum_x \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)} - \lambda \sum_a |H_t(a)| \frac{\partial}{\partial H_t(a)}$$

$$= \frac{\partial \left( \sum_x \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)} - \lambda \, \text{sign}(H_t(a)).$$

$$\underbrace{\qquad}_{I} \qquad \underbrace{\qquad}_{II}$$

$$\mathbf{I}$$

$$= \frac{\partial \left( \bar{\underset{x}{\sum}} \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)}$$

$$= H_t(a) + a\left( \bar{\underset{x}{\sum}} (q_*(x) - \beta_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \right)$$

$$= H_t(a) + a\left( \bar{\underset{x}{\sum}} (\pi_t(x)(q_*(x) - \beta_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \right)$$
$$\overline{\pi_t(x)}$$

$$= H_t(a) + a \, \mathbb{E}\left[ (Q_t - \bar{Q}_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} \, \middle/ \, \pi_t(A_t) \right]$$

$$= H_t(a) + a \left[ (Q_t - \bar{Q}_t) \left( \underset{a=x}{\mathbb{1}} - \pi_t(a) \right) - \lambda \, \text{sign} \, (H_t(a)) \right]$$

$$\text{if} \quad \underset{a=x}{\mathbb{1}} = 1, \quad \text{then} \quad a = x.$$

$$H_{t+1}(A_t) = H_t(A_t) + a(Q_t - \bar{Q}_t)(1 - \pi_t(A_t)) - \lambda \, \text{sign} \, (H_t(a))$$

$$\text{else} = 0.$$

$$H_{t+1}(a) = H_t(a) + a(Q_t - \bar{Q}_t)(-\pi_t(a)) - \lambda \, \text{sign} \, (H_t(a))$$
$$= H_t(a) - a(Q_t - \bar{Q}_t)(\pi_t(a) - \lambda \, \text{sign} \, (H_t(a)).$$

$$\text{5b. add } \ell_2.$$

$$H_{t+1}(a) = H_t(a) + a \frac{\partial \mathbb{E}[Q_t]}{\partial H_t(a)}.$$

$$\frac{\partial \mathbb{E}[Q_t]}{\partial H_t(a)} = \frac{\partial \left( \bar{\underset{x}{\sum}} \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)} - \lambda \left( \bar{\underset{x}{\sum}} [H_t(a)] \frac{\partial}{\partial H_t(a)} \right)^2$$

$$= \frac{\partial \left( \bar{\underset{x}{\sum}} \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)} - 2\lambda \, \text{sign}(H_t(a))$$
$$\underline{\mathbf{I}} \qquad\qquad\qquad \underline{\mathbf{II}}.$$

$$\underline{\mathbf{I}}$$

$$= \frac{\partial \left( \bar{\bar{\underset{x}{\sum}}} \pi_t(x) \cdot q_*(x) \right)}{\partial H_t(a)}.$$

$$= H_t(a) + a\left( \bar{\underset{x}{\sum}} (q_*(x) - \beta_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \right).$$

$$= H_t(a) + a\left( \bar{\underset{x}{\sum}} (\pi_t(x)(q_*(x) - \beta_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \right)$$
$$\overline{\pi_t(x)}.$$

$$= H_t(a) + a \, \mathbb{E}\left[ (Q_t - \bar{Q}_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} \, \middle/ \, \pi_t(A_t) \right].$$

$$= H_t(a) + a\left[(R_t - \bar{R_t})\left(\mathbb{1}_{a=t} - \pi_t(a)\right) - 2\lambda\, \text{sign}(H_t(a))\right].$$

if $\mathbb{1}_{a=x} = 1$, then $a = x$.

$$H_{t+1}(A_t) = H_t(A_t) + a(R_t - \bar{R_t})(1 - \pi_t(A_t)) - 2\lambda\, \text{sign}(H_t(a)).$$

else $= 0$.

$$H_{t+1}(a) = H_t(a) + a(R_t - \bar{R_t}) * (-\pi_t(a)) - 2\lambda\, \text{sign}(H_t(a))$$
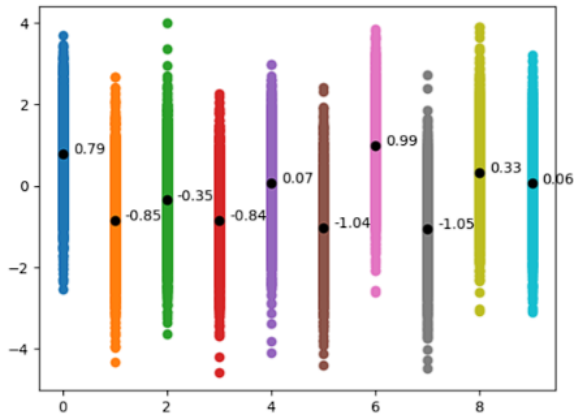$$= H_t(a) - a(R_t - \bar{R_t})(\pi_t(a)) - 2\lambda\, \text{sign}(H_t(a)).$$

## 6.

```
[72]: bandits = np.arange(0,10)
```

```
[73]: mean = 0
      std = 1
      nArms = 10 #n number of bandits
      start_time = time.time()
      iterations = 2000 # number of iterations
      plays = 1000 #number of plays per iterations


      rewards = np.random.normal(mean, std, nArms)
```

```
[74]: ran_rewards = np.array([ np.random.normal(rewards,std,iterations) for rewards in rewards]) #2000 interations is used
```

```
[75]: for index in bandits :
          plt.scatter(np.full(iterations,index),ran_rewards[index])
          plt.text(index+0.2,rewards[index],str(round(rewards[index],2)))
      plt.plot(bandits, rewards,'o', color='black')
      plt.show()
```

## 6a)

"A. Implement Greedy Method"

```
[79]: running_reward_sum = np.copy(reward_estimates) # _per_action_per_bandit

      rewards = []
      #rewards.append(0)
      rewards.append(np.mean(initial_reward_estimates)) # step 1

      epsilon = 0.0 # Greedy method
      for plays in range(2,plays):
          sum_of_reward = 0
          for problem_index in range(iterations):
              if np.random.random() > epsilon: # Greedy Selection
                  maxval = np.amax(reward_estimates[problem_index])
                  maxval_indices = np.ravel(np.array(np.where(reward_estimates[problem_index] == maxval)))
                  random_choice = np.random.choice( maxval_indices ) # breaking ties randomly
              else :
                  random_choice = np.random.randint(nArms)

              #print(str(problem_index),str(step),str(random_choice))
              running_reward_sum[problem_index][random_choice] += np.random.normal(testbed[problem_index][random_choice],1)
              action_count[problem_index][random_choice] += 1
              avg_reward = running_reward_sum[problem_index][random_choice] / action_count[problem_index][random_choice]
              reward_estimates[problem_index][random_choice] = avg_reward

              sum_of_reward += avg_reward

          rewards.append((sum_of_reward)/iterations)
```
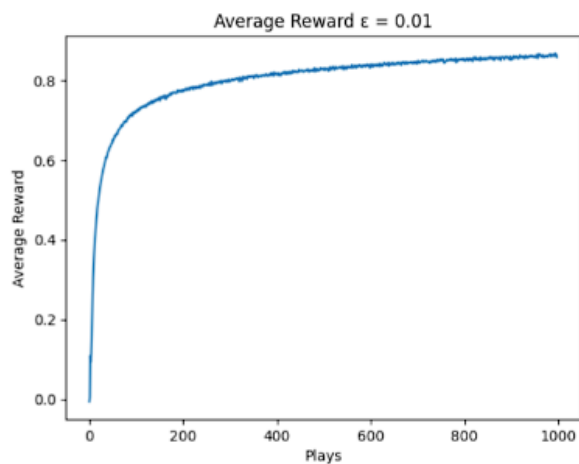
```
[80]: plt.plot(np.arange(plays),rewards)
```

[80]: [<matplotlib.lines.Line2D at 0x28f0cf18730>]

# 6b.

## B.Implement Epsilon-Greedy Method

```
[ ]:
```

```
[111]: running_reward_sum = np.copy(reward_estimates) # _per_action_per_bandit

       rewards = []
       #rewards.append(0)
       rewards.append(np.mean(initial_reward_estimates)) # step 1

       epsilon = 0.01 # Greedy method
       for plays in range(2,plays):
           sum_of_reward = 0
           for problem_index in range(iterations):
               if np.random.random() > epsilon: # Greedy Selection
                   maxval = np.amax(reward_estimates[problem_index])
                   maxval_indices = np.ravel(np.array(np.where(reward_estimates[problem_index] == maxval)))
                   random_choice = np.random.choice( maxval_indices ) # breaking ties randomly
               else :
                   random_choice = np.random.randint(nArms)

               #print(str(problem_index),str(step),str(random_choice))
               running_reward_sum[problem_index][random_choice] += np.random.normal(testbed[problem_index][random_choice],1)
               action_count[problem_index][random_choice] += 1
               avg_reward = running_reward_sum[problem_index][random_choice] / action_count[problem_index][random_choice]
               reward_estimates[problem_index][random_choice] = avg_reward

               sum_of_reward += avg_reward

           rewards.append((sum_of_reward)/iterations)
```
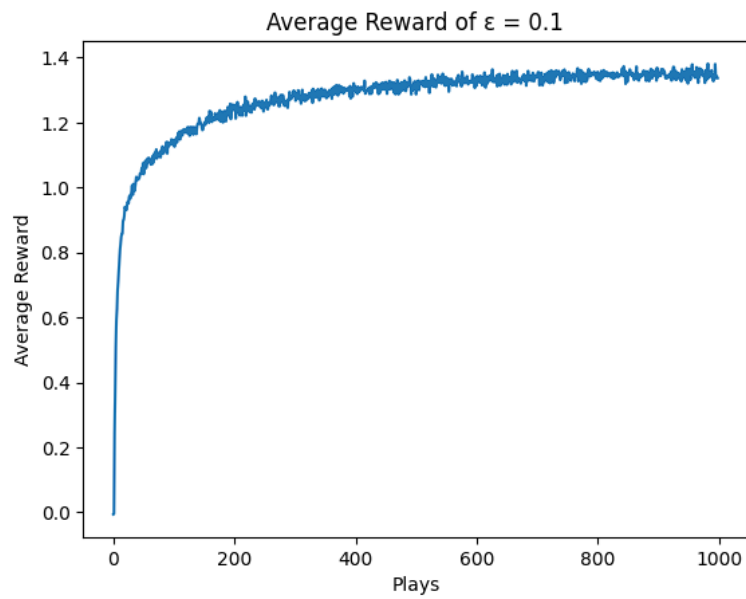
```
[112]: plt.xlabel('Plays')
       plt.ylabel('Average Reward')
       plt.title('Average Reward ε = 0.01')
       plt.plot(np.arange(plays),rewards, label=f"Epsilon = {epsilon}")
```

```
[112]: [<matplotlib.lines.Line2D at 0x28f0d59a9a0>]
```

Average Reward of ε = 0.1

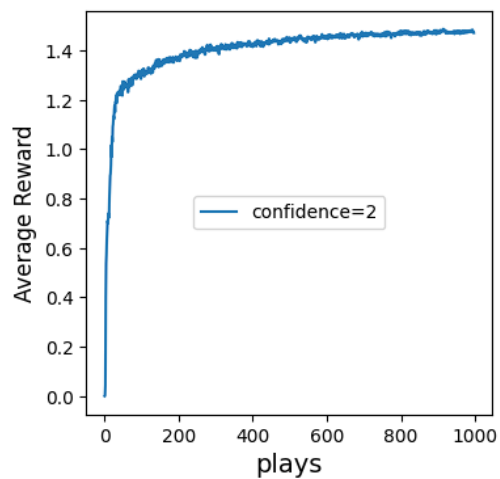## 6c.
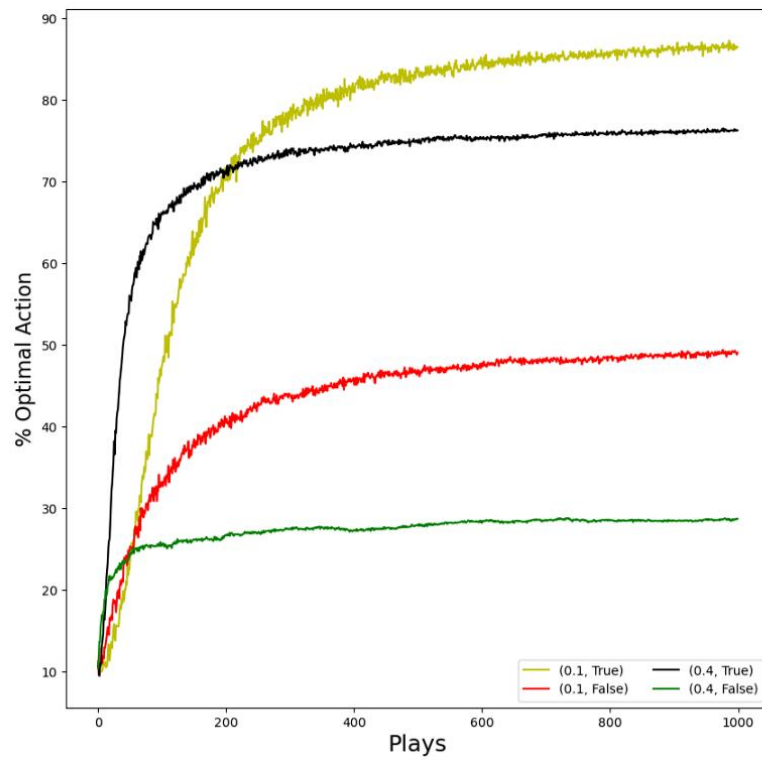
```
[126]:  ucb_result = []
        ucb_c(plays, iterations, nArms, 0 , 2,  testbed, initial_reward_estimates, ucb_result)

[130]:  plt.figure(figsize=(4,4))
        plt.xlabel('plays', fontsize=14)
        plt.ylabel('Average Reward', fontsize=12)
        plt.plot(np.arange(plays), ucb_result[0], label=("confidence=2"))
        plt.legend(loc='center', ncol=2)
        plt.show()
```
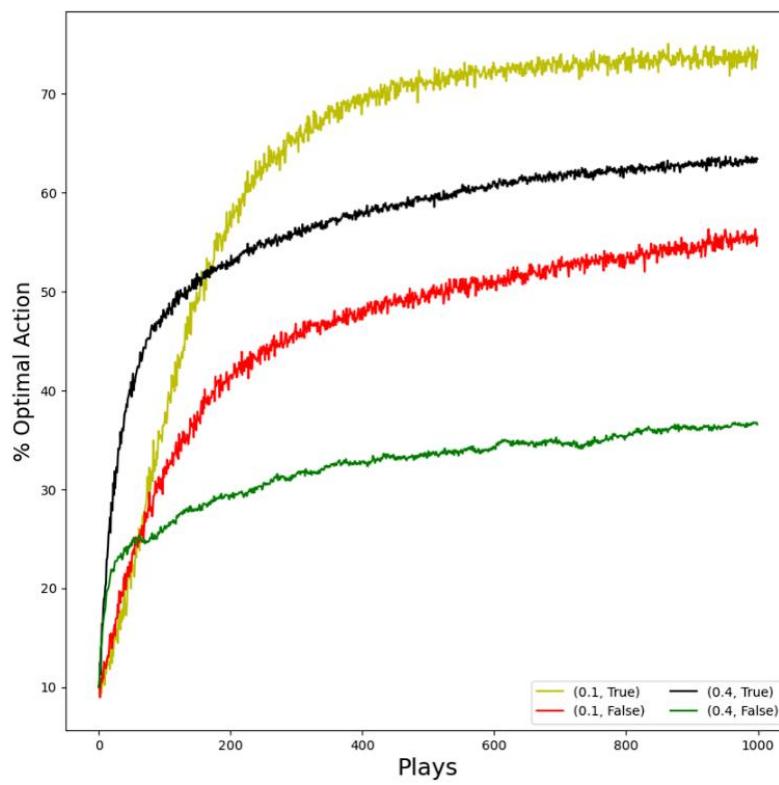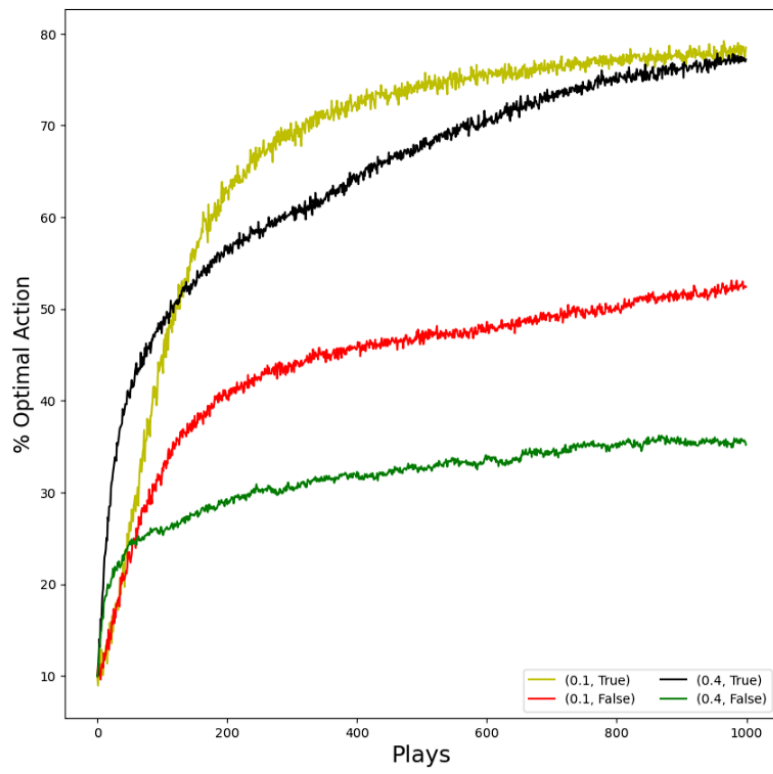
**6d.**



**6e.**

**6f.**

**7.** $b_t = p \, b + a_{t+1}$   $a_1 = -1$, $a_2 = 2$, $a_3 = 6$, $a_4 = 3$, $a_5 = 2$.

$b_5 = 0$   $\qquad T = 5$, $y = 0,5$

$b_4 = 0,5 \cdot b_5 + a_5$   $\qquad b_3 = 0,5 \cdot b_4 + a_4$   $\qquad b_2 = 0,5 \cdot b_3 + a_3$

$\quad = (0,5 \cdot 0) + 2$   $\qquad \quad = (0,5 \cdot 2) + 3$   $\qquad \quad = (0,5 \cdot 3) + 6$

$\quad = 2.$   $\qquad\qquad\qquad = 4$   $\qquad\qquad\qquad = 7,5$

$b_1 = 0,5 \cdot b_2 + a_2$   $\qquad\qquad b_0 = 0,5 \cdot b_1 + a_1$

$\quad = (0,5 \cdot 7,5) + 2 = 5,75$   $\qquad = (0,5 \cdot 5,75) + (-1)$

$\qquad\qquad\qquad\qquad\qquad\qquad = 1,875$

---

**8.** $y = 0,9$   $a_1 = 2$. Infinite Number of terms : $b_t = \sum_{h=0}^{\infty} y^h = \dfrac{1}{1-y}$,

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $1_s = 1$ time step.

$b_t = \sum_{h=0}^{\infty} y^h = \dfrac{1}{1-0,9}$

$\quad = \dfrac{1}{0,1}$  $(7_s)$   $\qquad\qquad b_1 = \dfrac{7}{0,1} + 7 = 77$

$\quad = \dfrac{7}{0,1}$   $\qquad\qquad\qquad b_0 = \dfrac{7}{0,1} + 2 = 72.$

---

**9.** $b_t = \sum_{u=0}^{\infty} y^h \cdot \dfrac{1}{1-y}$   $\qquad b_t = a_{t+1} + p\,a_{t+2} + y^2 a_{t+3} + \ldots$ ①

$\qquad\qquad\qquad\qquad\qquad y\,b_t = y\,a_{t+1} + y^2 a_{t+2} + y^3 a_{t+3} + \ldots y^{t+1} a_{t+n}.$

$\qquad\qquad\qquad\qquad \text{if } a = r.$

$\qquad\qquad\qquad\qquad y\,b_t = y\,r + y^2 r + y^3 r + \ldots y^{h+1} r$  ②

$b_t - y\,b_t = r - y^{h+1} r$

$b_t(1-y) = r(1 - y^{h+1})$

$\qquad b_t = \dfrac{r(1-y^{h+1})}{(1-y)}$

if $h \to \infty$   $b_t = \dfrac{r}{1-y}$   $\qquad a_t = r = 1$.

$\qquad\qquad\quad = \dfrac{1}{1-y}.$

## 10. b



## 10. d

11. Reward Hypothesis can be describe as maximizing the satisfaction from interacting a certain action in the environment. Reward can be obtain from every action as a feedback to the agent. The goal from this hypothesis is to maximizing the accumulated reward, althagh the current or immediate reward may be small. The agent must balance between short term benefits and long terms objectives. Agent must nadigate this trade-off, which often prioritizing actions that may yield smaller immediate reward, but lead to longer cumulative returns or reward in the future.