



**EXPERIMENT ASSESSMENT**

ACADEMIC YEAR 2023-24

**Course:** Reinforcement Learning Lab

**Course code:** CSDOL8013

**Year:** BE **Sem:** VIII

<b>Experiment No.:</b> 2
<b>Aim:</b> - Implementing a multi-armed bandit problem using UCB.
<b>Name:</b>
<b>Roll Number:</b>
<b>Date of Performance:</b>
<b>Date of Submission:</b>

**Evaluation**

Performance Indicator	Max. Marks	Marks Obtained
Performance	5	
Understanding	5	
Journal work and timely submission.	10	
<b>Total</b>	<b>20</b>	

Performance Indicator	Exceed Expectations (EE)	Meet Expectations (ME)	Below Expectations (BE)
Performance	5	3	2
Understanding	5	3	2
Journal work and timely submission.	10	8	4

**Checked by**

**Name of Faculty** : Rujuta Vartak

**Signature** :

**Date** :



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Science & Engineering (Data Science)

---

### EXPERIMENT 2

**Aim:** Implementing a multi-armed bandit problem using UCB.

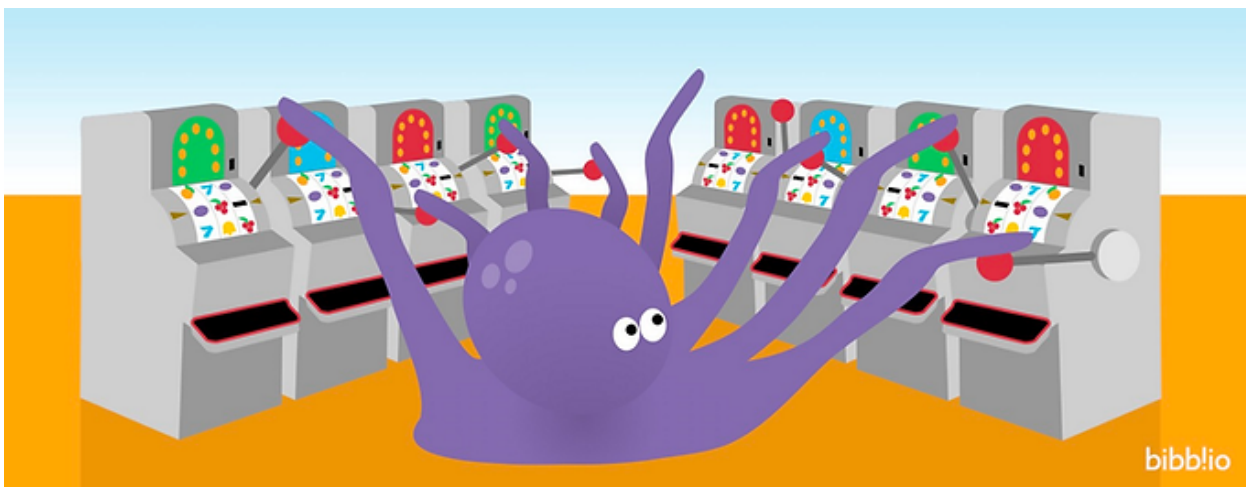
**Objective:** objective of UCB is to learn how the algorithm balances exploration by choosing arms with uncertain outcomes (based on confidence intervals) and exploitation by favoring arms with high estimated rewards.

#### **Theory:**

The Multi-Armed Bandit Problem (MAB) is a fundamental problem in the field of reinforcement learning and decision-making under uncertainty. The problem involves a gambler who has to choose among several slot machines (also called "one-armed bandits") with unknown payout probabilities. The gambler's objective is to maximize his or her total payout by choosing the best slot machine to play.

The MAB problem is commonly used in various fields, such as clinical trials, online advertising, and recommendation systems. The MAB problem can be seen as a trade-off between exploration and exploitation. Exploration refers to the gambler trying out different slot machines to learn about their payout probabilities, while exploitation refers to the gambler sticking to the slot machine that has shown the highest payout probability so far.

A bandit is defined as someone who steals your money. A one-armed bandit is a simple slot machine wherein you insert a coin into the machine, pull a lever, and get an immediate reward. But why is it called a bandit? It turns out all casinos configure these slot machines in such a way that all gamblers end up losing money!



A multi-armed bandit is a complicated slot machine wherein instead of 1, there are several levers that a gambler can pull, with each lever giving a different return. The probability distribution for the reward corresponding to each lever is different and is unknown to the gambler.



### Upper Confidence Bounds(UCB)

Upper Confidence Bounds (UCB) is a popular algorithm for solving the multi-armed bandit problem. It is similar to the  $\epsilon$ -greedy algorithm but uses a different strategy for balancing exploration and exploitation.

In UCB, the algorithm selects the arm with the highest upper confidence bound at each time step. The upper confidence bound is a measure of how uncertain the algorithm is about the reward distribution for that arm. It is calculated as follows:

**UCB = estimated mean reward + exploration bonus**

The exploration bonus is a function of the number of times the arm has been selected and the total number of times all arms have been selected. It is designed to encourage exploration of arms that have been selected less frequently, while still selecting arms with high estimated rewards.

The UCB algorithm starts by selecting each arm once to establish initial estimates of the reward distribution. It then selects the arm with the highest upper confidence bound at each subsequent time step. As more data is collected, the estimates of the reward distribution become more accurate and the algorithm becomes more confident in its arm selection.

With UCB, ' $A_t$ ', the action chosen at time step ' $t$ ', is given by:

$$A_t = \underset{a}{\operatorname{argmax}} \left( Q_t(a) + c \sqrt{\frac{\ln(t)}{N_t(a)}} \right)$$

Exploit      Explore

where;

- $Q_t(a)$  is the estimated value of action ' $a$ ' at time step ' $t$ '.
- $N_t(a)$  is the number of times that action ' $a$ ' has been selected, prior to time ' $t$ '.
- ' $c$ ' is a confidence value that controls the level of exploration.

One advantage of UCB over the  $\epsilon$ -greedy algorithm is that it does not require a parameter to specify the level of exploration. The algorithm automatically adjusts the level of exploration based on the uncertainty of the reward distribution for each arm.

In summary, Upper Confidence Bounds (UCB) is an algorithm for solving the multi-armed bandit problem

that balances exploration and exploitation by selecting the arm with the highest upper confidence bound at each time step. It is an effective algorithm for solving the problem and does not require a parameter to specify the level of exploration.



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Science & Engineering (Data Science)

---

### Code:

```
# importing essential libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
# importing our dataset "ads.csv"
dataset=pd.read_csv("/content/ads.csv")
# viewing the first five observations of our dataset
dataset.head(5)
import random
N=10000                                     # Number of records in our
dataset                                    dataset
d=10                                       # Number of versions
ads_selected=[]
total_reward=0
for n in range(0,N):
    ad=random.randrange(d)
    ads_selected.append(ad)
    reward=dataset.values[n,ad]
    total_reward=total_reward + reward
plt.hist(ads_selected)
plt.title('Histogram of ads selections')
plt.xlabel('Ads')
plt.ylabel('Number of times each ad was selected')
plt.show()
# Interpretation:-
# So from the above histogram plot for random implementation method we can
see that:
# 1] from the random selection method all the ads are getting multiple random
selections.
# 2] so we are not able to properly find the ad which is selected maximum
number of times.
# 3] Also in our histogram we can't properly see , which ad is being selected
maximum number of times.
# 4] And as a result we cant find which ad to display in future.
# 5] so, from random selection method we can't achieve our goal of obtaining
maximum selected ad, which is the biggest drawback
# of random selection method.
# implementing UCB
import math
```



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Science & Engineering (Data Science)

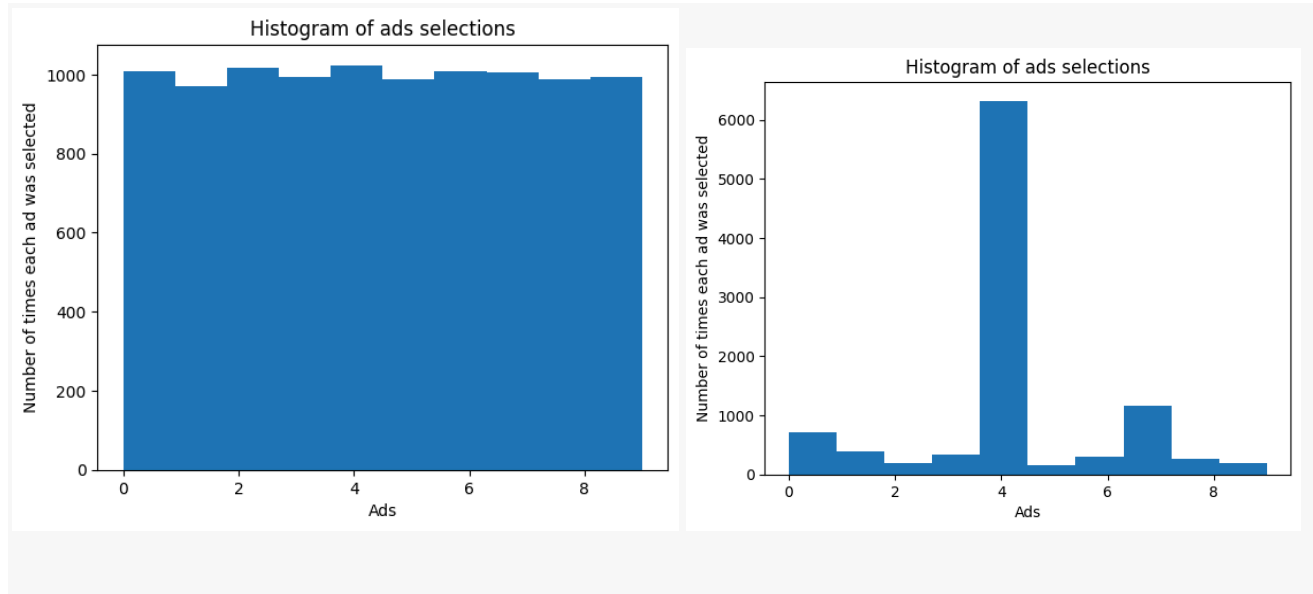
```
N=10000                                # Number of records in our
dataset                                dataset
d=10                                   # Number of versions
ads_selected=[]
number_of_selections=[0]*d             # Ni(n)
sums_of_rewards=[0]*d                 # Ri(n)
total_reward=0
for n in range(0,N):
    ad=0
    max_upper_bound=0
    for i in range(0,d):
        if (number_of_selections[i]>0):
            average_reward=sums_of_rewards[i]/number_of_selections[i]
# Ri /Ni
            delta_i=math.sqrt(3/2*math.log(n+1)/number_of_selections[i])
# confidence interval
            upper_bound=average_reward + delta_i
        else:
            upper_bound=1e400
        if upper_bound > max_upper_bound:
            max_upper_bound= upper_bound
        ad= i
# step-2 ends
ads_selected.append(ad)
# step-3 starts
number_of_selections[ad]=number_of_selections[ad]+1
reward=dataset.values[n,ad]
sums_of_rewards[ad]=sums_of_rewards[ad] + reward
total_reward= total_reward + reward
# we got the maximum values asked in step-3
print(ads_selected)
plt.hist(ads_selected)
plt.title('Histogram of ads selections')
plt.xlabel('Ads')
plt.ylabel('Number of times each ad was selected')
plt.show()
```



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Science & Engineering (Data Science)

### Output:



### Conclusion:

#### 1. Explain Example of Multi-arm Bandit Problem.

The multi-armed bandit problem metaphorically reflects a scenario where a gambler faces multiple slot machines (arms) in a casino, each with unknown payout probabilities. The gambler's objective is to maximize their total reward by selecting which machines to play and how frequently, with only partial information about the machines' reward distributions. The challenge lies in striking a balance between exploring different machines to learn about their potential rewards (exploration) and exploiting the knowledge gained to maximize immediate rewards (exploitation).

#### 2. Compare $\epsilon$ -greedy and UCB.

Both  $\epsilon$ -greedy and UCB algorithms address the exploration-exploitation trade-off in the multi-armed bandit problem, UCB offers a more principled approach by incorporating uncertainty into the decision-making process, leading to more efficient and adaptive exploration strategies.