

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

- (1) State the Doubling trick formally. Explain what problem it solves and show how. (5 points)

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Q.1) Doubling Trick Algorithm:—
Algorithm follows below steps,

1. For $m=1, 2, \dots$
2. Run a new instance of algorithm on the 2^m rounds $t \leq 2^m$, i.e. $(2^{m+1} - 1)$ with optimal learning rate. (for an algorithm that runs for 2^m steps)

Problem Doubling Trick Solved:—
For randomized weighted majority algorithm zero regret bound can be achieved if $\epsilon^* = \sqrt{\frac{\ln d}{T}}$; where, $T \rightarrow$ No. of rounds
 $d \rightarrow$ No. of experts

Mathematically it states that— for weighted majority algorithm, if $T \rightarrow \infty$ then $R_A(T) \rightarrow 0$
 \rightarrow No. of rounds \rightarrow Algorithm regret

It actually means that— as the number of rounds increase our algorithm can track the best expert & achieve the zero regret.

But— the problem here is, in real life scenario it is impossible to reach $T \rightarrow \infty$ for any practical case.

As the number of rounds in practical scenario is always finite.

To overcome this issue one approach is to come up with some finite Horizon Regret Bound.

To simulate the algorithm for $T \rightarrow \infty$ we can perform below task,

- a) Let the algorithm fix T as some finite value, say N .
- b) Complete the algorithm for first N rounds & keep track of the loss suffered.
- c) Restart the algorithm again & run it for same N no. of rounds.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

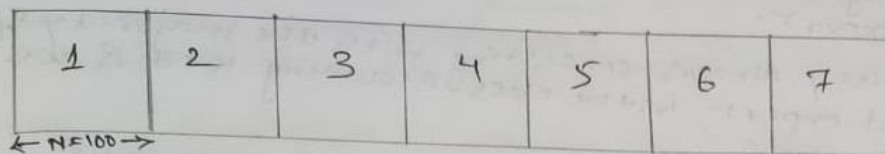
Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Below is the 7 iteration & each having 100 round as per our example,



Loss suffered in each round is same.

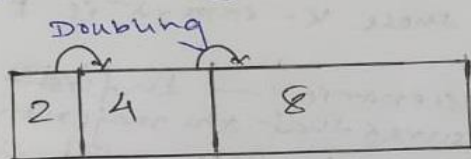
In above diagram we have chosen $N=100$ & restarted the algorithm in each 100 round. We have observed the same loss has been suffered as at each block algorithm is restarting.

Problem in this set up is we will never achieve good (improvement) in the suffered loss.

That also means if we make $T \rightarrow \infty$, in this way there will be no improvement in regret.

To solve this issue, "Doubling Trick" algorithm came into picture.

If we change the settings in such a way that if we change the no. of iteration round in every steps that means if doubled, diagram will be as follows,



if we use this settings then total regret will be,

$$\sum_{i=0}^{\log T} \sqrt{2^i \log d}$$

[d \rightarrow No. of expert
 $T \rightarrow$ no. of round]

$$= \sqrt{\log d} \sum_{i=0}^{\log T} 2^{i/2}$$

$$\therefore \left[\text{As per formula we know, } \sum_{i=1}^K a^i = \frac{a^{K+1} - a}{a - 1} \right]$$

Simplifying the eqn we got, $= \sqrt{\log d} \cdot \sqrt{T}$

which indicate that loss has suffered in order of $\sqrt{T \log d}$.

$$\text{loss suffered} \approx O(\sqrt{T \log d})$$

Conclusion:— So, in this settings we will modify the ϵ value at every round after doubling the round & we will get the regret guarantee, which is in the order of $O(\sqrt{T \log d})$. This is called "Doubling Trick".

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

- (2) Consider the usual online learning protocol with a finite set of d experts. Consider an adversary who is always consistent with the majority of k out of d experts in the class (assume $k < d$ and k is odd). What is the worst case number of mistakes for any algorithm for this problem? Can you think of an algorithm which achieves the same? (10 points)

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Q.2)

According to the problem statement - there are a finite set of d experts.

Adversary always consisting with the majority of k out of d experts, in the class assuming $k < d$ & k is an odd number.

Main goal of an algorithm (learner) is to make prediction that are as accurate as the best possible expert in hindsight.

As per the problem statement, we can think this problem in broadly 2 way,

a) Either k experts are consistent with adversary.

b) $\frac{k+1}{2}$ experts are consistent with adversary.

If we explain this with some example, let's consider,

if $k=11$

→ 11 experts consistent with adversary

→ 6 experts consistent with adversary

If we consider both cases mentioned above,

a) If we go with majority vote algorithm:-

	d	k	our decision	Truth	Conclusion
divided by 2	64	11			
→ 32			0	1	[Mistake]
→ 16			0	1	[Mistake]

In both the steps, there is a mistake happen, so we are half the experts in each steps are removed.

Since, we know that at least k are consistent, so we have to wait till they get the majority vote.

Now for the explained example, there are 2 mistakes are made.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

similarly,

d	K	NO. of mistakes
64	11	2
64	1	6

Looking at that table we can observe that, as the value of the K increases, no. of mistakes decrease. Since as per problem statement we have to find out the worst number of mistakes, K value has to be decrease. actually K should be minimum.

b) As per the problem statement we are open to debate on any algorithm for worst case mistakes, if we are start picking our experts randomly & adversary force us to make all the output as incorrect then we will end up with making $(d-K)$ mistakes.

And if we consider the worst case $K=1$ then that will become $(d-1)$ mistakes.

We know,

$$(d-1) \text{ mistake} > \log_2(d)$$

We can say that if we do random guesses for our output then we can achieve the worst mistakes of $(d-1)$.

Now if we consider the next interpretation,

$$\frac{K+1}{2} \text{ experts are consistent with adversary}$$

→ Then, from the previous conclusion we can draw the fact that,

$$d - \left(\frac{K+1}{2}\right) \text{ no. of mistakes will happen for worst case scenario}$$

$$i.e. (2d - K - 1 / 2)$$

So, in this scenario K can have minimum value of 3 such that worst case mistakes can be $(d-2)$.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

(3) Assume that you want to estimate the probability of a coin coming up heads. Let the true probability be μ .

- Assume $\mu = 0.75$. If you want to make a statement with at least 99% confidence that the estimate $\hat{\mu}$ satisfies $|\mu - \hat{\mu}| \leq 0.1$, how many times would you have to toss the coin? (4 points)
- If the true μ changed from 0.75 to 0.5, how does your answer from the first part change? (2 points)
- If the confidence changed from 99% to 90%, how does your answer from the first part change? (2 points)
- If the tolerance change from 0.1 to 0.01, how does your answer from the first part change? (2 points)

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Q3) Assuming that I want to estimate the probability of a coin coming up as heads.
Let the true probability be μ .
 $\hat{\mu}$ is the estimates which satisfy $|\mu - \hat{\mu}| \leq 0.1$
According to the "Hoeffding's Inequality" provides an upper bound on the probability that the sum of bounded independent random variables deviate from its expected value by more than a certain amount.

$$P(|\mu - \hat{\mu}| > \epsilon) \leq 2 \exp(-2n\epsilon^2) \quad \dots (i)$$

Let consider

$$P(|\mu - \hat{\mu}| > \epsilon) = \alpha$$

then equation (i) can be written as,

$$\alpha \leq 2 \exp(-2n\epsilon^2)$$

take log in both side of the equation,

$$\Rightarrow \ln \alpha \leq \ln [2 \exp(-2n\epsilon^2)]$$

$$\Rightarrow \ln \alpha \leq \ln 2 + (-2n\epsilon^2)$$

$$\Rightarrow \ln \alpha \leq \ln 2 - 2n\epsilon^2$$

$$\Rightarrow n \geq \frac{\ln(2/\alpha)}{2\epsilon^2} \quad \dots (ii)$$

PART 1:-

Now according to the problem statement,

for 99% confidence bound $\alpha = \frac{(1-0.99)}{2} = 0.005$

$$\& \epsilon = 0.1$$

plotting the value in equation (ii)

$$n \geq \frac{\ln(2/0.005)}{2 \times (0.1)^2}$$

$$\Rightarrow n \geq \frac{5.29832}{0.02}$$

$$\Rightarrow n \geq 264.9 \approx 265$$

\therefore we have to toss the coin for at least 265 times.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

PART 2:-

If true μ changes from 0.75 to 0.5 then there will not be having any impact.

Because if we consider "Hoeffding Inequality" then it is clear that n value does not depend on the true probability.

$$n \geq \frac{\ln(2/\alpha)}{2\epsilon^2} \quad ; \text{ no true probability term exists in this equation.}$$

PART 3:-

As per the equation (ii) from the general equation derived using "Hoeffding Inequality", we know that,

$$n \geq \frac{\ln(2/\alpha)}{2\epsilon^2}$$

As the confidence level changes from 99% to 90% we can place,

$$\alpha = 0.1 \quad \& \quad \epsilon = 0.1$$

So, putting this value in the formula we got,

$$n \geq \frac{\ln(2/\alpha)}{2\epsilon^2}$$
$$\Rightarrow n \geq \frac{\ln(2/0.1)}{2 \times (0.1)^2}$$

$$\Rightarrow n \geq \frac{2.9957}{0.02} \approx 149.785 \approx 150$$

\therefore In this settings we have to toss the coin for at least 150 times.

PART 4:-

In this settings tolerance change from 0.1 to 0.01. Then in this case values will be like as follows,

$$\alpha = 0.01 \quad \& \quad \epsilon = 0.01$$

$$n \geq \frac{\ln(2/\alpha)}{2\epsilon^2}$$

$$\Rightarrow n \geq \frac{\ln(2/0.01)}{2 \times (0.01)^2}$$

$$\Rightarrow n \geq 26491.586 \approx 26492$$

\therefore So in this settings we have to flip the coin at least 26492 times.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

- (4) An *epsilon*-greedy strategy for the stochastic multi-armed bandits set up exploits the current best arm with probability $(1-\epsilon)$ and explores with a small probability ϵ . Consider a problem instance with 10 arms where the reward for the i -th ($i = 1, \dots, 10$) arm is Beta distributed with parameters $\alpha_i = 5, \beta_i = 5 * i$.
- Implement the *epsilon*-greedy algorithm and compare it with the performance of the UCB and the EXP-3 algorithm. (5 points)
 - Comment on your observations about the regret plots obtained in the previous part. (2.5 points)
 - If you vary ϵ , how does the regret change? (2.5 points)

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Screenshot:

Part 1 - Implement the epsilon-greedy algorithm and compare it with the performance of the UCB and the EXP-3 algorithm. (5 points)

Importing Library

```
1 import numpy as np
2 import matplotlib as mat
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import pandas as pd
```

Common Variable

```
1 # Setting seed to control randomness for reproducibility
2 seed = 48
3
4 # Initializing random number generator with the defined seed
5 np.random.seed(seed)
6
7 # Defining common variables for running the algorithm
8 num_arms = 10 # Number of arms/options in the multi-armed bandit problem
9 num_iterations = 1000 # Number of iterations to run the algorithm
10
11 # Parameters of the Beta distributions for each arm's reward
12 # Success parameters (alpha) for each arm's Beta distribution, all set to 5
13 success_params = np.ones(num_arms) * 5
14 # Failure parameters (beta) for each arm's Beta distribution, increasing by 5 for each arm
15 failure_params = np.arange(5, 5*(num_arms+1), 5)
16
17 # Reshape failure_params to match the shape of success_params
18 # Ensuring the failure_params array has same shape as success_params
19 failure_params = np.resize(failure_params, success_params.shape) |
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
1 # Epsilon greedy Algorithm
2 def epsilon_greedy_policy(epsilon_val, num_options, num_iterations, success_params, failure_params):
3
4     """
5     Implements the epsilon-greedy algorithm for a multi-armed bandit problem.
6
7     Args:
8         epsilon_val (float): The probability of choosing a random option (exploration).
9         num_options (int): The number of options, or "arms," of the bandit.
10        num_iterations (int): The number of iterations to run the algorithm.
11        success_params (array_like): Parameters of the Beta distribution for success for each option.
12        failure_params (array_like): Parameters of the Beta distribution for failure for each option.
13
14    Returns:
15        np.array: The cumulative regret after each time step.
16
17    This function operates over a specified number of iterations. At each time step,
18    it chooses the option with the highest average reward with probability (1-epsilon),
19    and chooses a random option with probability epsilon. The function calculates the
20    "regret" of each option as the difference in success probability from the optimal
21    option. The function returns the cumulative regret after each time step.
22    """
23
24    # Initialize cumulative rewards and number of times each option is chosen
25    total_rewards = np.zeros(num_options)
26    option_uses = np.zeros(num_options)
27
28    # Initialize rewards and regret lists for tracking over time
29    reward_tracker = []
30    regret_tracker = []
31
32    # Determine the "best" option under current understanding
33    optimal_option = np.argmax((success_params / (success_params + failure_params)))
34
35    # Begin the main loop for the specified number of time steps
36    for time_step in range(num_iterations):
37        # Choose the option with highest average reward most of the time (probability 1-epsilon)
38        if np.random.random() < (1 - epsilon_val):
39            # Add small value to prevent division by zero
40            selected_option = np.argmax(total_rewards / (option_uses + 1e-6))
41        else: # but choose randomly some of the time (probability epsilon)
42            selected_option = np.random.randint(num_options)
43
44        # Generate reward based on a Beta distribution (specific to each option)
45        current_reward = np.random.beta(success_params[selected_option], failure_params[selected_option])
46
47        # Update rewards and usage for the selected option
48        total_rewards[selected_option] += current_reward
49        option_uses[selected_option] += 1
50
51        # Record current reward and regret (difference in success probability from optimal option)
52        reward_tracker.append(current_reward)
53        regret_tracker.append(abs(success_params[optimal_option] / (success_params[optimal_option]
54                                                                    + failure_params[optimal_option])
55                               - success_params[selected_option] / (success_params[selected_option]
56                                                                    + failure_params[selected_option])))
57
58    # Return cumulative regret over time
59    return np.cumsum(regret_tracker)
```


ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
62 # UCB algorithm
63 def UCB_algorithm(num_options, num_iterations, success_params, failure_params):
64     """
65     Implements the Upper Confidence Bound (UCB) algorithm for a multi-armed bandit problem.
66
67     Args:
68         num_options (int): The number of options, or "arms," of the bandit.
69         num_iterations (int): The number of iterations to run the algorithm.
70         success_params (array_like): Parameters of the Beta distribution for success for each option.
71         failure_params (array_like): Parameters of the Beta distribution for failure for each option.
72
73     Returns:
74         np.array: The cumulative regret after each time step.
75
76     The function operates over a specified number of iterations. At each time step,
77     it calculates the UCB value for each option, and chooses the option with the
78     highest UCB value. The UCB value for an option is its average reward plus an
79     uncertainty term, which decreases as that option is chosen more frequently.
80
81     The function calculates the "regret" of each option as the difference in success
82     probability from the optimal option. The function returns the cumulative regret
83     after each time step.
84     """
85
86     # Initialize cumulative rewards and number of times each option is chosen
87     total_rewards = np.zeros(num_options)
88     option_uses = np.ones(num_options)
89
90     # Initialize rewards and regret lists for tracking over time
91     reward_list = []
92     regret_list = []
93
94     # Determine the "best" option under current understanding
95     optimal_option = np.argmax((success_params / (success_params + failure_params)))
96
97     # Begin the main loop for the specified number of time steps, starting from num_options
98     # (as each option needs to be tried once)
99     for time_step in range(num_options, num_iterations):
100         # Calculate UCB values for each option
101         ucb_values = total_rewards / option_uses + np.sqrt(2 * np.log(time_step) / option_uses)
102         # Choose the option that currently has the highest UCB value
103         selected_option = np.argmax(ucb_values)
104
105         # Generate reward based on a Beta distribution (specific to each option)
106         current_reward = np.random.beta(success_params[selected_option], failure_params[selected_option])
107
108         # Update rewards and usage for the selected option
109         total_rewards[selected_option] += current_reward
110         option_uses[selected_option] += 1
111
112         # Record current reward and regret (difference in success probability from optimal option)
113         reward_list.append(current_reward)
114         regret_list.append(abs(success_params[optimal_option] / (success_params[optimal_option]
115                                                                + failure_params[optimal_option])
116                             - success_params[selected_option] / (success_params[selected_option]
117                                                                + failure_params[selected_option])))
118
119     # Return cumulative regret over time
120     return np.cumsum(regret_list)
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
122 # EXP-3 algorithm
123 def EXP3_algorithm(num_options, num_iterations, success_params, failure_params, learning_rate):
124     """
125     Implements the Exponential-weight algorithm for Exploration and Exploitation (EXP3) for a multi-armed bandit problem.
126
127     Args:
128         num_options (int): The number of options, or "arms," of the bandit.
129         num_iterations (int): The number of iterations to run the algorithm.
130         success_params (array_like): Parameters of the Beta distribution for success for each option.
131         failure_params (array_like): Parameters of the Beta distribution for failure for each option.
132         learning_rate (float): Learning rate parameter for the EXP3 algorithm.
133
134     Returns:
135         np.array: The cumulative regret after each time step.
136
137     The function operates over a specified number of iterations. At each time step,
138     it calculates the probabilities for each option based on their weights and selects
139     an option randomly according to the calculated probabilities.
140
141     The function calculates the "regret" of each option as the difference in success
142     probability from the optimal option. The function returns the cumulative regret
143     after each time step.
144     """
145     # Initialize cumulative rewards and weights for each option
146     total_rewards = np.zeros(num_options)
147     option_weights = np.ones(num_options)
148
149     # Initialize rewards and regret lists for tracking over time
150     reward_records = []
151     regret_records = []
152
153     # Determine the "best" option under current understanding
154     optimal_option = np.argmax((success_params / (success_params + failure_params)))
155
156     # Begin the main loop for the specified number of time steps
157     for time_step in range(num_iterations):
158         # Calculate probabilities for each option
159         option_probs = (1 - learning_rate) * (option_weights / np.sum(option_weights)) + learning_rate / num_options
160         # Select an option randomly according to the calculated probabilities
161         selected_option = np.random.choice(num_options, p=option_probs)
162
163         # Generate reward based on a Beta distribution (specific to each option)
164         current_reward = np.random.beta(success_params[selected_option], failure_params[selected_option])
165
166         # Update the total reward for the selected option
167         total_rewards[selected_option] += current_reward
168
169         # Estimate the reward and update the weight of the selected option
170         estimated_reward = total_rewards[selected_option] / option_probs[selected_option]
171         option_weights[selected_option] *= np.exp(learning_rate * estimated_reward / num_options)
172
173         # Normalize weights to prevent underflow
174         option_weights /= np.sum(option_weights)
175
176         # Record current reward and regret (difference in success probability from optimal option)
177         reward_records.append(current_reward)
178         regret_records.append(abs(success_params[optimal_option] / (success_params[optimal_option] + failure_params[optimal_option]) -
179                                success_params[selected_option] / (success_params[selected_option] + failure_params[selected_option])))
180
181     # Return cumulative regret over time
182     return np.cumsum(regret_records)
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
84 def plot_cumulative_regret(epsilon_greedy_cumulative_regret, ucb_cumulative_regret, exp3_cumulative_regret, epsilon_values,
85 """
86     Plots the cumulative regret for Epsilon-Greedy, UCB, and EXP3 algorithms.
87
88     Args:
89         epsilon_greedy_cumulative_regret (array_like): Cumulative regret values from the Epsilon-Greedy algorithm.
90         ucb_cumulative_regret (array_like): Cumulative regret values from the UCB algorithm.
91         exp3_cumulative_regret (array_like): Cumulative regret values from the EXP3 algorithm.
92         epsilon_values (float): Epsilon value used in the Epsilon-Greedy algorithm.
93         learning_params (float): Learning parameter used in the EXP3 algorithm.
94
95     The function creates a line plot for each set of regret values, using different colors
96     for each algorithm. It sets labels for the x-axis and y-axis, a title for the plot, and
97     a legend. It also includes grid lines for easier viewing of the data points.
98     """
99     # Create a colormap using the 'cool' colormap from the matplotlib library
100     cmap = mat.colormaps['cool']
101     # Generate three colors from the colormap. These colors are evenly spaced
102     # from the beginning (0) to the end (1) of the colormap.
103     colors = cmap(np.linspace(0, 1, 3))
104     # Set the figure size for the plot
105     plt.figure(figsize=(10,7))
106     # Plot the cumulative regret for the Epsilon-Greedy algorithm.
107     # The color for this line is the first color from the colormap.
108     plt.plot(epsilon_greedy_cumulative_regret, label="Epsilon-Greedy", color=colors[0])
109     # Plot the cumulative regret for the UCB algorithm.
110     # The color for this line is the second color from the colormap.
111     plt.plot(ucb_cumulative_regret, label="UCB", color=colors[1])
112     # Plot the cumulative regret for the EXP-3 algorithm.
113     # The color for this line is the third color from the colormap.
114     plt.plot(exp3_cumulative_regret, label="EXP-3", color=colors[2])
115     # Set the label for the x-axis
116     plt.xlabel("No. of Rounds", fontsize=14)
117     # Set the label for the y-axis
118     plt.ylabel("Cumulative Regret", fontsize=14)
119     # Set the title for the plot
120     plt.title(f"Comparison of Epsilon-Greedy, UCB, and EXP-3 algorithms for epsilon value {epsilon_values} & learning param")
121     # Display the legend
122     plt.legend(fontsize=12)
123     # Display grid lines
124     plt.grid(True, linestyle='--', alpha=0.6)
125     # Show the plot
126     plt.show()
127
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
def process_results(df, results, algorithm_name, epsilon_val, learning_param):
    """
    Process the results of a multi-armed bandit algorithm, recording data for each set of parameters.

    Args:
        df (pandas.DataFrame): The DataFrame to which the results will be appended.
        results (list): A list of tuples containing the parameters and regret values for each run of the algorithm.
        algorithm_name (str): The name of the algorithm used.

    Each tuple in `results` should contain three elements:
    1. epsilon_val (float): The epsilon value used in the algorithm run.
    2. learning_param (float): The learning parameter used in the algorithm run.
    3. regrets (list): A list of regret values from the algorithm run.

    The function identifies the maximum and minimum regret in each run and their respective iterations,
    creates a new DataFrame with these details, and appends it to the provided DataFrame `df`.

    Returns:
        df_analyser (pandas.DataFrame): The original DataFrame, updated with the new data.
    """
    # Iterate through each tuple in the results list
    for epsilon_val, learning_param, regrets in results:
        # Find the maximum regret value and its corresponding iteration number
        max_regret_iter, max_regret = max(enumerate(regrets), key=lambda x: x[1])
        # Find the minimum regret value and its corresponding iteration number
        min_regret_iter, min_regret = min(enumerate(regrets), key=lambda x: x[1])

        # Create a new DataFrame to store the results of the current tuple
        new_data = pd.DataFrame({
            'Algorithm': [algorithm_name], # The algorithm name
            'Epsilon': [epsilon_val], # The epsilon value used
            'Learning_Rate': [learning_param], # The Learning parameter used
            'Max_Regret': [max_regret], # The maximum regret value
            'Max_Regret_Iter': [max_regret_iter + 1], # The iteration number for the maximum regret value
            'Min_Regret': [min_regret], # The minimum regret value
            'Min_Regret_Iter': [min_regret_iter + 1] # The iteration number for the minimum regret value
        })

        # Append the new data to the existing DataFrame
        df_analyser = pd.concat([df, new_data], ignore_index=True)

    # Return the updated DataFrame
    return df_analyser
```


ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

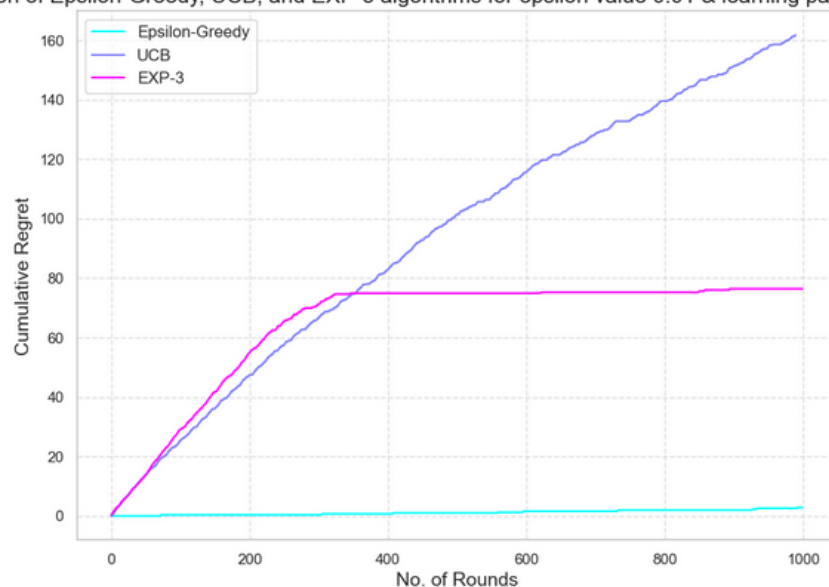
Prepared By: Aloy Banerjee

Roll Number: CH22M503

Execution of the algorithm

```
1 # The probability used in the epsilon-greedy algorithm to determine
2 # whether to explore (choose a random arm) or exploit (choose the best arm so far).
3 epsilon_values = [0.01]
4 learning_params = [0.01]
5
6 # Prepare empty lists to store results
7 epsilon_greedy_cumulative_regret, ucb_cumulative_regret, exp3_cumulative_regret = [], [], []
8
9 # Prepare empty lists to store results
10 epsilon_greedy_results, ucb_results, exp3_results = [], [], []
11
12
13 # Run simulations for each combination of epsilon and learning parameters
14 for epsilon_val in epsilon_values:
15     for learning_param in learning_params:
16         epsilon_greedy_cumulative_regret = epsilon_greedy_policy(epsilon_val, num_arms, num_iterations, success_params, fai
17         ucb_cumulative_regret = UCB_algorithm(num_arms, num_iterations, success_params, failure_params)
18         exp3_cumulative_regret = EXP3_algorithm(num_arms, num_iterations, success_params, failure_params, learning_param)
19         plot_cumulative_regret(epsilon_greedy_cumulative_regret, ucb_cumulative_regret, exp3_cumulative_regret, epsilon_val
20         epsilon_greedy_results.append((epsilon_val, learning_param, epsilon_greedy_cumulative_regret))
21         ucb_results.append((epsilon_val, learning_param, ucb_cumulative_regret))
22         exp3_results.append((epsilon_val, learning_param, exp3_cumulative_regret))
23
```

Comparison of Epsilon-Greedy, UCB, and EXP-3 algorithms for epsilon value 0.01 & learning parameter value 0.01



ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Loading the data in dataframe for analysis

```
1 # Create an empty DataFrame
2 df_Analyser = pd.DataFrame(columns=['Algorithm', 'Epsilon', 'Learning_Rate', 'Max_Regret',
3                                     'Max_Regret_Iter', 'Min_Regret', 'Min_Regret_Iter'])
4
5 # Process results for each algorithm and add to the DataFrame
6 df_Analyser = process_results(df_Analyser, epsilon_greedy_results, 'Epsilon-Greedy', epsilon_values[0],
7                               learning_params[0])
8 df_Analyser = process_results(df_Analyser, ucb_results, 'UCB', epsilon_values[0], learning_params[0])
9 df_Analyser = process_results(df_Analyser, exp3_results, 'Exp 3', epsilon_values[0], learning_params[0])
10 df_Analyser_sorted = df_Analyser.sort_values(['Epsilon', 'Learning_Rate'], ascending=[True, True])
```

Viewed the sorted dataframe containing the information about the regret for different algorithm for selected Epsilon and Learning rate

```
1 df_Analyser_sorted
```

87]:

	Algorithm	Epsilon	Learning_Rate	Max_Regret	Max_Regret_Iter	Min_Regret	Min_Regret_Iter
0	Epsilon-Greedy	0.01	0.01	2.842083	993	0.0	1
1	UCB	0.01	0.01	161.713203	990	0.0	1
2	Exp 3	0.01	0.01	76.372691	895	0.3	1

Part 2 : Comment on your observations about the regret plots obtained in the previous part. (2.5 points)

Conclusion of Part 1 :

- Looking at the plot for three algorithm reveals that the value of cumulative regret over the number of rounds,
 - for UCB : Remain in increasing trend.
 - for Exp-3 : Initially the regret value is increasing but eventually it flattened out after certain round.
 - for Epsilon Greedy : On the contrary cumulative regret remain very close to zero and change very slowly.
- We are getting same indication from the table displayed in the last cell also.

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Part 3 : Execution of the algorithm for different values of ϵ

```
1 seed = 20
2 np.random.seed(seed)
3 num_iterations = 1000
4 # Different epsilon values to explore
5 epsilon_values_to_explore = [0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5,
6                             0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95]
7
8 # List to store the cumulative regrets for each epsilon value
9 cumulative_regrets_for_each_epsilon = []
10
11 # Run the epsilon-greedy policy for each epsilon value and store the cumulative regrets
12 for epsilon in epsilon_values_to_explore:
13     cumulative_regret = epsilon_greedy_policy(epsilon, num_arms, num_iterations, success_params, failure_params)
14     cumulative_regrets_for_each_epsilon.append(cumulative_regret)
15
16 # Set the style and context of the plot to make it look more attractive
17 sns.set(style='whitegrid', context='notebook')
18
19 # Set the color palette for the plot
20 colors = [mat.colormaps['tab20'](i) for i in range(20)]
21
22 # Create a larger figure
23 plt.figure(figsize=(12, 8))
24
25 # Loop over epsilon values
26 for i, epsilon in enumerate(epsilon_values_to_explore):
27     # Plot the cumulative regret for each epsilon, with a different color and line style for each one
28     plt.plot(cumulative_regrets_for_each_epsilon[i], color=colors[i], linestyle='--', label=f"Epsilon = {epsilon}")
29
30 # Set labels and title with larger fonts
31 plt.xlabel("No. of Rounds", fontsize=14)
32 plt.ylabel("Cumulative Regret", fontsize=14)
33 plt.title("Epsilon-Greedy Algorithm with Regret and Various Epsilon Values", fontsize=16)
34
35 # Display the legend and increase its font size
36 plt.legend(fontsize=12)
37
38 # Add a slight grid for better readability
39 plt.grid(True, linestyle='--', alpha=0.6)
40
41 # Show the plot
42 plt.show()
```

ID6002W: Online Learning and Reinforcement Learning

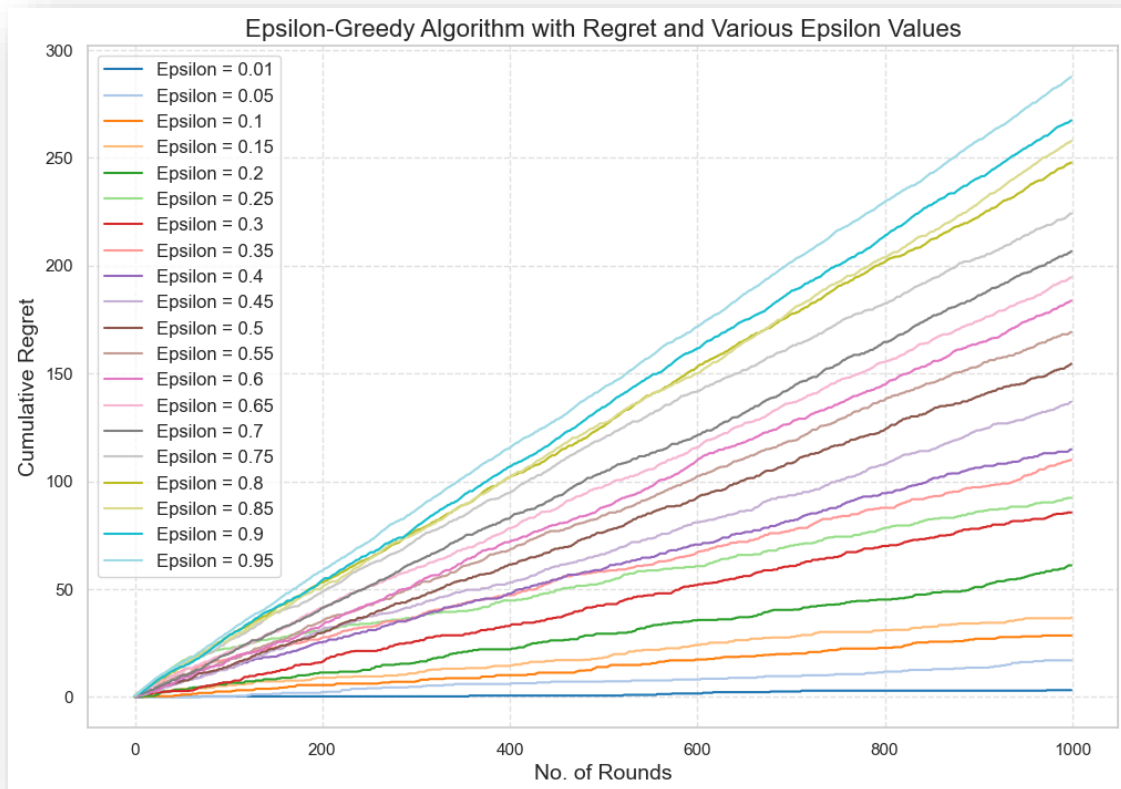
Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503



Conclusion Part 3 :

- Various epsilon values (0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95) demonstrate various regret patterns in the Epsilon-Greedy algorithm. In terms of minimizing regrets and observing that the regret is growing its slope as and when we are raising the epsilon value over the number of rounds, smaller epsilon values (for example, 0.01) tend to perform better than bigger epsilon values (for example, 0.95). This suggests that a more experimental strategy (lower epsilon) can eventually result in greater overall performance.

ID6002W: Online Learning and Reinforcement Learning
Mid Term Exam
Course Instructor: Arun Rajkumar.
Exam Date: 18-06-2023
Prepared By: Aloy Banerjee
Roll Number: CH22M503

Code Text:

Part 1 - Implement the epsilon-greedy algorithm and compare it with the performance of the UCB and the EXP-3 algorithm. (5 points)

Importing Library

```
import numpy as np
import matplotlib as mat
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
```

Common Variable

```
# Setting seed to control randomness for reproducibility
seed = 48

# Initializing random number generator with the defined seed
np.random.seed(seed)

# Defining common variables for running the algorithm
num_arms = 10 # Number of arms/options in the multi-armed bandit problem
num_iterations = 1000 # Number of iterations to run the algorithm

# Parameters of the Beta distributions for each arm's reward
# Success parameters (alpha) for each arm's Beta distribution, all set to 5
success_params = np.ones(num_arms) * 5
# Failure parameters (beta) for each arm's Beta distribution, increasing by 5 for each arm
failure_params = np.arange(5, 5*(num_arms+1), 5)

# Reshape failure_params to match the shape of success_params
# Ensuring the failure_params array has same shape as success_params
failure_params = np.resize(failure_params, success_params.shape)
```

Common Method

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
# Epsilon greedy Algorithm
```

```
def epsilon_greedy_policy(epsilon_val, num_options, num_iterations,  
success_params, failure_params):
```

```
    """  
    Implements the epsilon-greedy algorithm for a multi-armed bandit  
    problem.
```

```
    Args:
```

```
        epsilon_val (float): The probability of choosing a random option  
        (exploration).
```

```
        num_options (int): The number of options, or "arms," of the  
        bandit.
```

```
        num_iterations (int): The number of iterations to run the  
        algorithm.
```

```
        success_params (array_like): Parameters of the Beta distribution  
        for success for each option.
```

```
        failure_params (array_like): Parameters of the Beta distribution  
        for failure for each option.
```

```
    Returns:
```

```
        np.array: The cumulative regret after each time step.
```

```
    This function operates over a specified number of iterations. At each  
    time step,
```

```
        it chooses the option with the highest average reward with probability  
        (1-epsilon),
```

```
        and chooses a random option with probability epsilon. The function  
        calculates the
```

```
        "regret" of each option as the difference in success probability from  
        the optimal
```

```
        option. The function returns the cumulative regret after each time  
        step.
```

```
    """
```

```
    # Initialize cumulative rewards and number of times each option is  
    chosen
```

```
    total_rewards = np.zeros(num_options)
```

```
    option_uses = np.zeros(num_options)
```

```
    # Initialize rewards and regret lists for tracking over time
```

```
    reward_tracker = []
```

```
    regret_tracker = []
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
# Determine the "best" option under current understanding
optimal_option = np.argmax((success_params / (success_params +
failure_params)))

# Begin the main loop for the specified number of time steps
for time_step in range(num_iterations):
    # Choose the option with highest average reward most of the time
    (probability 1-epsilon)
    if np.random.random() < (1 - epsilon_val):
        # Add small value to prevent division by zero
        selected_option = np.argmax(total_rewards / (option_uses + 1e-
6))
    else: # but choose randomly some of the time (probability
epsilon)
        selected_option = np.random.randint(num_options)

    # Generate reward based on a Beta distribution (specific to each
option)
    current_reward = np.random.beta(success_params[selected_option],
failure_params[selected_option])

    # Update rewards and usage for the selected option
    total_rewards[selected_option] += current_reward
    option_uses[selected_option] += 1

    # Record current reward and regret (difference in success
probability from optimal option)
    reward_tracker.append(current_reward)
    regret_tracker.append(abs(success_params[optimal_option] /
(success_params[optimal_option] + failure_params[optimal_option]) -
success_params[selected_option] / (success_params[selected_option] +
failure_params[selected_option])))

# Return cumulative regret over time
return np.cumsum(regret_tracker)

# UCB algorithm
def UCB_algorithm(num_options, num_iterations, success_params,
failure_params):
    """
    Implements the Upper Confidence Bound (UCB) algorithm for a multi-
armed bandit problem.
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Args:

num_options (int): The number of options, or "arms," of the bandit.
num_iterations (int): The number of iterations to run the algorithm.
success_params (array_like): Parameters of the Beta distribution for success for each option.
failure_params (array_like): Parameters of the Beta distribution for failure for each option.

Returns:

np.array: The cumulative regret after each time step.

The function operates over a specified number of iterations. At each time step,

it calculates the UCB value for each option, and chooses the option with the

highest UCB value. The UCB value for an option is its average reward plus an

uncertainty term, which decreases as that option is chosen more frequently.

The function calculates the "regret" of each option as the difference in success

probability from the optimal option. The function returns the cumulative regret

after each time step.

"""

Initialize cumulative rewards and number of times each option is chosen

total_rewards = np.zeros(num_options)

option_uses = np.ones(num_options)

Initialize rewards and regret lists for tracking over time

reward_list = []

regret_list = []

Determine the "best" option under current understanding

optimal_option = np.argmax((success_params / (success_params + failure_params)))

Begin the main loop for the specified number of time steps, starting from num_options

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
# (as each option needs to be tried once)
for time_step in range(num_options, num_iterations):
    # Calculate UCB values for each option
    ucb_values = total_rewards / option_uses + np.sqrt(2 *
np.log(time_step) / option_uses)
    # Choose the option that currently has the highest UCB value
    selected_option = np.argmax(ucb_values)

    # Generate reward based on a Beta distribution (specific to each
option)
    current_reward = np.random.beta(success_params[selected_option],
failure_params[selected_option])

    # Update rewards and usage for the selected option
    total_rewards[selected_option] += current_reward
    option_uses[selected_option] += 1

    # Record current reward and regret (difference in success
probability from optimal option)
    reward_list.append(current_reward)
    regret_list.append(abs(success_params[optimal_option] /
(success_params[optimal_option] + failure_params[optimal_option]) -
success_params[selected_option] / (success_params[selected_option] +
failure_params[selected_option])))

    # Return cumulative regret over time
    return np.cumsum(regret_list)

# EXP-3 algorithm
def EXP3_algorithm(num_options, num_iterations, success_params,
failure_params, learning_rate):
    """
    Implements the Exponential-weight algorithm for Exploration and
Exploitation (EXP3) for a multi-armed bandit problem.

    Args:
        num_options (int): The number of options, or "arms," of the
bandit.
        num_iterations (int): The number of iterations to run the
algorithm.
        success_params (array_like): Parameters of the Beta distribution
for success for each option.
        failure_params (array_like): Parameters of the Beta distribution
for failure for each option.
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

learning_rate (float): Learning rate parameter for the EXP3 algorithm.

Returns:

np.array: The cumulative regret after each time step.

The function operates over a specified number of iterations. At each time step, it calculates the probabilities for each option based on their weights and selects an option randomly according to the calculated probabilities.

The function calculates the "regret" of each option as the difference in success probability from the optimal option. The function returns the cumulative regret after each time step.

```
"""
# Initialize cumulative rewards and weights for each option
total_rewards = np.zeros(num_options)
option_weights = np.ones(num_options)

# Initialize rewards and regret lists for tracking over time
reward_records = []
regret_records = []

# Determine the "best" option under current understanding
optimal_option = np.argmax((success_params / (success_params +
failure_params)))

# Begin the main loop for the specified number of time steps
for time_step in range(num_iterations):
    # Calculate probabilities for each option
    option_probs = (1 - learning_rate) * (option_weights /
np.sum(option_weights)) + learning_rate / num_options
    # Select an option randomly according to the calculated
probabilities
    selected_option = np.random.choice(num_options, p=option_probs)

    # Generate reward based on a Beta distribution (specific to each
option)
    current_reward = np.random.beta(success_params[selected_option],
failure_params[selected_option])
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
# Update the total reward for the selected option
total_rewards[selected_option] += current_reward

# Estimate the reward and update the weight of the selected option
estimated_reward = total_rewards[selected_option] /
option_probs[selected_option]
option_weights[selected_option] *= np.exp(learning_rate *
estimated_reward / num_options)

# Normalize weights to prevent underflow
option_weights /= np.sum(option_weights)

# Record current reward and regret (difference in success
probability from optimal option)
reward_records.append(current_reward)
regret_records.append(abs(success_params[optimal_option] /
(success_params[optimal_option] + failure_params[optimal_option]) -
success_params[selected_option] / (success_params[selected_option] +
failure_params[selected_option])))

# Return cumulative regret over time
return np.cumsum(regret_records)
```

```
def plot_cumulative_regret(epsilon_greedy_cumulative_regret,
ucb_cumulative_regret, exp3_cumulative_regret, epsilon_values,
learning_params):
```

```
    """
    Plots the cumulative regret for Epsilon-Greedy, UCB, and EXP3
    algorithms.
```

```
    Args:
```

```
        epsilon_greedy_cumulative_regret (array_like): Cumulative regret
        values from the Epsilon-Greedy algorithm.
```

```
        ucb_cumulative_regret (array_like): Cumulative regret values from
        the UCB algorithm.
```

```
        exp3_cumulative_regret (array_like): Cumulative regret values from
        the EXP3 algorithm.
```

```
        epsilon_values (float): Epsilon value used in the Epsilon-Greedy
        algorithm.
```

```
        learning_params (float): Learning parameter used in the EXP3
        algorithm.
```

```
    The function creates a line plot for each set of regret values, using
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
different colors
    for each algorithm. It sets labels for the x-axis and y-axis, a title
    for the plot, and
    a legend. It also includes grid lines for easier viewing of the data
    points.
    """
    # Create a colormap using the 'cool' colormap from the matplotlib
    library
    cmap = mat.colormaps['cool']
    # Generate three colors from the colormap. These colors are evenly
    spaced
    # from the beginning (0) to the end (1) of the colormap.
    colors = cmap(np.linspace(0, 1, 3))
    # Set the figure size for the plot
    plt.figure(figsize=(10,7))
    # Plot the cumulative regret for the Epsilon-Greedy algorithm.
    # The color for this line is the first color from the colormap.
    plt.plot(epsilon_greedy_cumulative_regret, label="Epsilon-Greedy",
    color=colors[0])
    # Plot the cumulative regret for the UCB algorithm.
    # The color for this line is the second color from the colormap.
    plt.plot(ucb_cumulative_regret, label="UCB", color=colors[1])
    # Plot the cumulative regret for the EXP-3 algorithm.
    # The color for this line is the third color from the colormap.
    plt.plot(exp3_cumulative_regret, label="EXP-3", color=colors[2])
    # Set the label for the x-axis
    plt.xlabel("No. of Rounds", fontsize=14)
    # Set the label for the y-axis
    plt.ylabel("Cumulative Regret", fontsize=14)
    # Set the title for the plot
    plt.title(f"Comparison of Epsilon-Greedy, UCB, and EXP-3 algorithms
    for epsilon value {epsilon_values} & learning parameter value
    {learning_params}", fontsize=16)
    # Display the Legend
    plt.legend(fontsize=12)
    # Display grid lines
    plt.grid(True, linestyle='--', alpha=0.6)
    # Show the plot
    plt.show()
```

```
def process_results(df, results, algorithm_name, epsilon_val,
    learning_param):
    """
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

Process the results of a multi-armed bandit algorithm, recording data for each set of parameters.

Args:

df (pandas.DataFrame): The DataFrame to which the results will be appended.

results (list): A list of tuples containing the parameters and regret values for each run of the algorithm.

algorithm_name (str): The name of the algorithm used.

Each tuple in `results` should contain three elements:

1. epsilon_val (float): The epsilon value used in the algorithm run.

2. learning_param (float): The learning parameter used in the algorithm run.

3. regrets (list): A list of regret values from the algorithm run.

The function identifies the maximum and minimum regret in each run and their respective iterations,

creates a new DataFrame with these details, and appends it to the provided DataFrame `df`.

Returns:

df_analyser (pandas.DataFrame): The original DataFrame, updated with the new data.

```
"""
# Iterate through each tuple in the results list
for epsilon_val, learning_param, regrets in results:
    # Find the maximum regret value and its corresponding iteration
    number
    max_regret_iter, max_regret = max(enumerate(regrets), key=lambda
x: x[1])
    # Find the minimum regret value and its corresponding iteration
    number
    min_regret_iter, min_regret = min(enumerate(regrets), key=lambda
x: x[1])

# Create a new DataFrame to store the results of the current tuple
new_data = pd.DataFrame({
    'Algorithm': [algorithm_name], # The algorithm name
    'Epsilon': [epsilon_val], # The epsilon value used
    'Learning_Rate': [learning_param], # The learning parameter
used
    'Max_Regret': [max_regret], # The maximum regret value
```


ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
        'Max_Regret_Iter': [max_regret_iter + 1], # The iteration
number for the maximum regret value
        'Min_Regret': [min_regret], # The minimum regret value
        'Min_Regret_Iter': [min_regret_iter + 1] # The iteration
number for the minimum regret value
    })

    # Append the new data to the existing DataFrame
    df_analyser = pd.concat([df, new_data], ignore_index=True)

    # Return the updated DataFrame
    return df_analyser
```

Execution of the algorithm

```
# The probability used in the epsilon-greedy algorithm to determine
# whether to explore (choose a random arm) or exploit (choose the best arm
so far).
epsilon_values = [0.01]
learning_params = [0.01]

# Prepare empty lists to store results
epsilon_greedy_cumulative_regret, ucb_cumulative_regret,
exp3_cumulative_regret = [], [], []

# Prepare empty lists to store results
epsilon_greedy_results, ucb_results, exp3_results = [], [], []

# Run simulations for each combination of epsilon and learning parameters
for epsilon_val in epsilon_values:
    for learning_param in learning_params:
        epsilon_greedy_cumulative_regret =
epsilon_greedy_policy(epsilon_val, num_arms, num_iterations,
success_params, failure_params)
        ucb_cumulative_regret = UCB_algorithm(num_arms, num_iterations,
success_params, failure_params)
        exp3_cumulative_regret = EXP3_algorithm(num_arms, num_iterations,
success_params, failure_params, learning_param)
        plot_cumulative_regret(epsilon_greedy_cumulative_regret,
ucb_cumulative_regret, exp3_cumulative_regret, epsilon_val,
learning_param)
        epsilon_greedy_results.append((epsilon_val, learning_param,
epsilon_greedy_cumulative_regret))
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

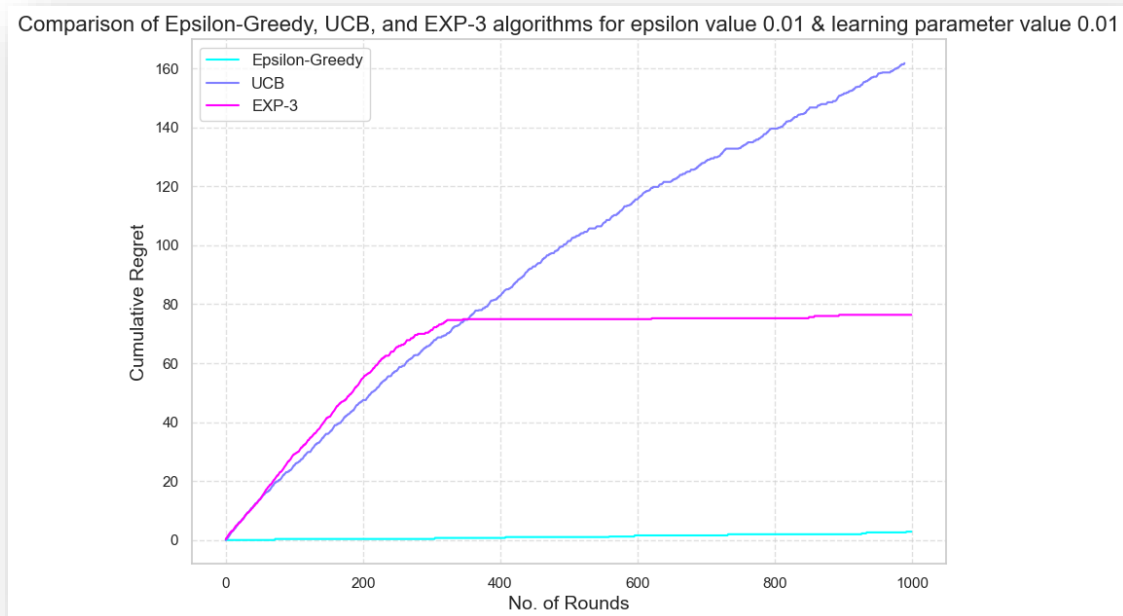
Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
ucb_results.append((epsilon_val, learning_param,
ucb_cumulative_regret))
exp3_results.append((epsilon_val, learning_param,
exp3_cumulative_regret))
```



Loading the data in dataframe for analysis

```
# Create an empty DataFrame
df_Analyser = pd.DataFrame(columns=['Algorithm', 'Epsilon',
'Learning_Rate', 'Max_Regret',
'Max_Regret_Iter', 'Min_Regret',
'Min_Regret_Iter'])

# Process results for each algorithm and add to the DataFrame
df_Analyser = process_results(df_Analyser, epsilon_greedy_results,
'Epsilon-Greedy', epsilon_values[0],
learning_params[0])
df_Analyser = process_results(df_Analyser, ucb_results, 'UCB',
epsilon_values[0], learning_params[0])
df_Analyser = process_results(df_Analyser, exp3_results, 'Exp 3',
epsilon_values[0], learning_params[0])
```

ID6002W: Online Learning and Reinforcement Learning

Mid Term Exam

Course Instructor: Arun Rajkumar.

Exam Date: 18-06-2023

Prepared By: Aloy Banerjee

Roll Number: CH22M503

```
df_Analyser_sorted = df_Analyser.sort_values(['Epsilon', 'Learning_Rate'],
ascending=[True, True])
```

Viewed the sorted dataframe containing the information about the regret for different algorithm for selected Epsilon and Learning rate

```
df_Analyser_sorted
```

	Algorithm	Epsilon	Learning_Rate	Max_Regret	Max_Regret_Iter \
0	Epsilon-Greedy	0.01	0.01	2.842063	993
1	UCB	0.01	0.01	161.713203	990
2	Exp 3	0.01	0.01	76.372691	895

	Min_Regret	Min_Regret_Iter
0	0.0	1
1	0.0	1
2	0.3	1

Part 2: Comment on your observations about the regret plots obtained in the previous part. (2.5 points)

Conclusion of Part 1:

- Looking at the plot for three algorithm reveals that the value of cumulative regret over the number of rounds,
 - for **UCB**: Remain in increasing trend.
 - for **Exp-3**: Initially the regret value increased but eventually it flattened out after a certain round.
 - for **Epsilon Greedy**: On the contrary cumulative regret remains very close to zero and changes very slowly.
- We are getting the same indication from the table displayed in the last cell also.

Part 3: Execution of the algorithm for different values of ϵ

```
# Different epsilon values to explore
epsilon_values_to_explore = [0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35,
0.4, 0.45, 0.5,
                                0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9,
0.95]
```

```
# List to store the cumulative regrets for each epsilon value
cumulative_regrets_for_each_epsilon = []
```

ID6002W: Online Learning and Reinforcement Learning
Mid Term Exam
Course Instructor: Arun Rajkumar.
Exam Date: 18-06-2023
Prepared By: Aloy Banerjee
Roll Number: CH22M503

```
# Run the epsilon-greedy policy for each epsilon value and store the cumulative regrets
for epsilon in epsilon_values_to_explore:
    cumulative_regret = epsilon_greedy_policy(epsilon, num_arms,
num_iterations, success_params, failure_params)
    cumulative_regrets_for_each_epsilon.append(cumulative_regret)

# Set the style and context of the plot to make it look more attractive
sns.set(style='whitegrid', context='notebook')

# Set the color palette for the plot
colors = [mat.colormaps['tab20'](iloop) for iloop in range(20)]

# Create a larger figure
plt.figure(figsize=(12, 8))

# Loop over epsilon values
for i, epsilon in enumerate(epsilon_values_to_explore):
    # Plot the cumulative regret for each epsilon, with a different color and line style for each one
    plt.plot(cumulative_regrets_for_each_epsilon[i], color=colors[i],
linestyle='--', label=f"Epsilon = {epsilon}")

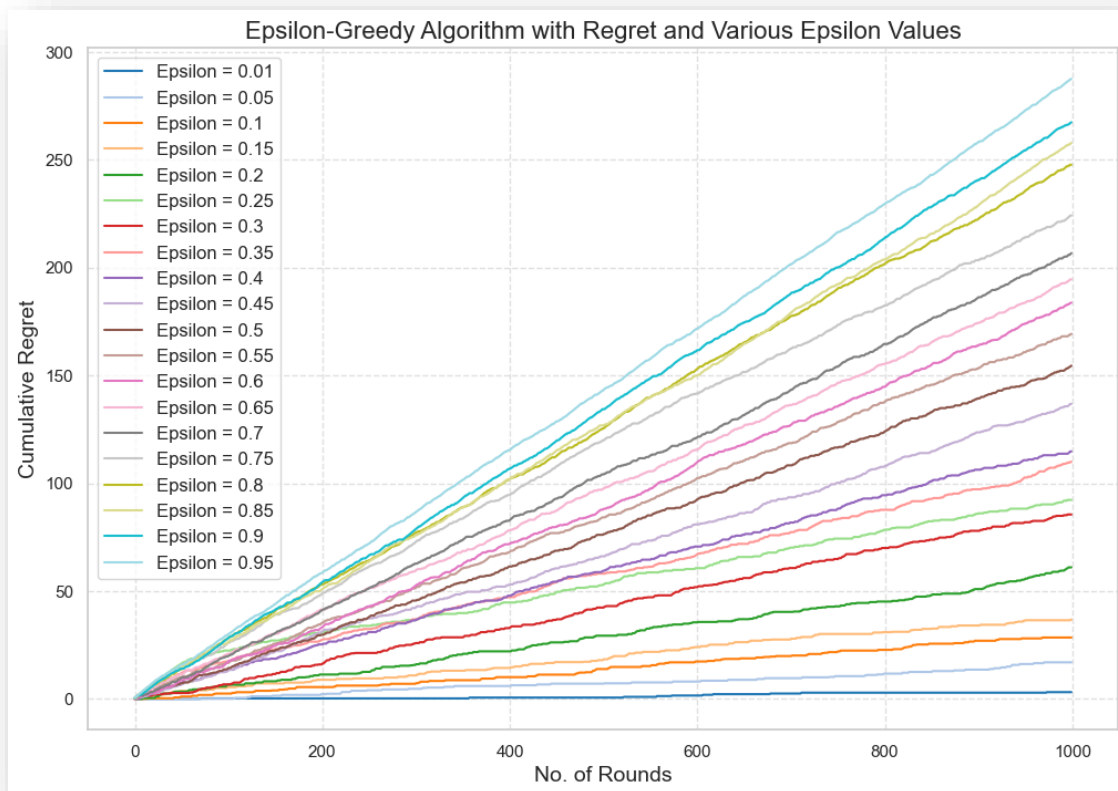
# Set Labels and title with larger fonts
plt.xlabel("No. of Rounds", fontsize=14)
plt.ylabel("Cumulative Regret", fontsize=14)
plt.title("Epsilon-Greedy Algorithm with Regret and Various Epsilon Values", fontsize=16)

# Display the Legend and increase its font size
plt.legend(fontsize=12)

# Add a slight grid for better readability
plt.grid(True, linestyle='--', alpha=0.6)

# Show the plot
plt.show()
```

ID6002W: Online Learning and Reinforcement Learning
Mid Term Exam
Course Instructor: Arun Rajkumar.
Exam Date: 18-06-2023
Prepared By: Aloy Banerjee
Roll Number: CH22M503



Conclusion Part 3:

- Various epsilon values (0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95) demonstrate various regret patterns in the Epsilon-Greedy algorithm. In terms of minimizing regrets and observing that the regret is growing its slope as and when we are raising the epsilon value over the number of rounds, smaller epsilon values (for example, 0.01) tend to perform better than bigger epsilon values (for example, 0.95). This suggests that a more experimental strategy (lower epsilon) can eventually result in greater overall performance.

End of Report