



**Subject: Reinforcement  
Learning**

**AY: 2023 – 24**

**Experiment 1  
Exploration Exploitation Dilemma**

**AIM:**

- a) To solve the exploration exploitation dilemma using epsilon greedy strategy
- b) To understand the effect of epsilon by comparing the different values of epsilon

**THEORY:**

With partial knowledge about future states and future rewards, our reinforcement learning agent will be in a dilemma on whether to exploit the partial knowledge to receive some rewards or it should explore unknown actions which could result in much larger rewards.

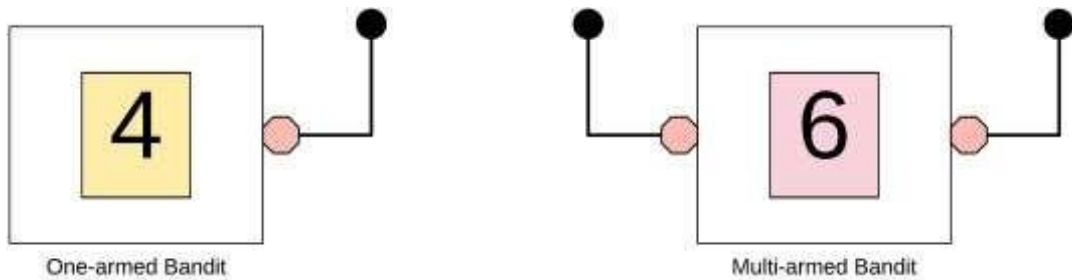
- Exploitation: Make the best decision given current information. The best long-term strategymay involve short-term sacrifices
- Exploration Gather more information. Gather enough information to make the best overalldecisions

However, we cannot choose to explore and exploit simultaneously. In order to overcome the Exploration-Exploitation Dilemma, we use the Epsilon Greedy Policy.

**MULTI ARMED BANDIT PROBLEM (MAB)**

In a multi-armed bandit problem (MAB) (or n-armed bandits), an Agent makes a choice from a set of actions. This choice results in a numeric reward from the Environment based on the selected action. In this specific case, the nature of the Environment is a stationary probability distribution. By stationary, we mean that the probability distribution is constant (or independent) across all states of the Environment. In other words, the probability distribution is unchanged as the state of the Environment changes. The goal of the Agent in a MAB problem is to maximize the rewards received from the Environment over a specified period.

The MAB problem is an extension of the “one-armed bandit” problem, which is represented as a slot machine in a casino. In the MAB setting, instead of a slot machine with one-lever, we have multi-levers. Each lever corresponds to an action the Agent can play. The goal of the Agent is to make plays that maximize its winnings (i.e., rewards) from the machine. The Agent will have to figure out the best levers (exploration) and then concentrate on the levers (exploitation) that will maximize its returns (i.e., the sum of the rewards).



Left: One-armed bandit. The slot machine has one lever that returns a numerical reward when played.

Right: Multi-armed bandits. The slot machine has multiple (n) arms, each returning a numerical reward when played. In a MAB problem, the reinforcement agent must balance exploration and exploitation to maximize returns.

### Epsilon-greedy

The agent does random exploration occasionally with probability  $\epsilon$  and takes the optimal action most of the time with probability  $1 - \epsilon$ .

Epsilon greedy method. At each step, a random number is generated by the model. If the number was lower than epsilon in that step (exploration area) the model chooses a random action and if it was higher than epsilon in that step (exploitation area) the model chooses an action based on what it learned.

Usually, epsilon is set to be around 10%. Epsilon-Greedy can be represented as follows:

$$A_t \leftarrow \begin{cases} \operatorname{argmax} Q_t(a) & \text{with probability } 1 - \epsilon \\ a \sim \operatorname{Uniform}(\{a_1 \dots a_k\}) & \text{with probability } \epsilon \end{cases}$$

The Action that the agent selects at time step  $t$ , will be a greedy action (exploit) with probability  $(1 - \epsilon)$  or may be a random action (explore) with probability of  $\epsilon$ .

**ALGORITHM:**

---

**Algorithm 2:** Epsilon-Greedy Action Selection

---

**Data:** Q: Q-table generated so far,  $\epsilon$ : a small number, S: current state

**Result:** Selected action

**Function** *SELECT-ACTION*(Q, S,  $\epsilon$ ) **is**

```
    n  $\leftarrow$  uniform random number between 0 and 1;  
    if  $n < \epsilon$  then  
        | A  $\leftarrow$  random action from the action space;  
    else  
        | A  $\leftarrow$  maxQ(S,.);  
    end  
    return selected action A;  
end
```

---

**LAB ASSIGNMENT TO DO:**

1. Create a multi armed bandit agent which would estimate the win rate using the epsilon-greedy strategy.
2. Understand the effect of the value of epsilon on the win rate by comparing the win rates corresponding to different values of epsilon. Hence draw conclusions.