

1.

State Space (S)

The state space for the autonomous drone delivery system can be defined with multiple dimensions to encapsulate its environment comprehensively. Possible components of the state space include:

- **Position:** The drone's current GPS coordinates (latitude, longitude, altitude).
- **Velocity:** The drone's speed and direction in three-dimensional space (X, Y, Z axes).
- **Battery Level:** The remaining battery percentage, crucial for flight endurance and operation.
- **Obstacle Information:** Sensor data outlining the distance and type of nearby obstacles such as buildings, trees, and other drones.
- **Destination Information:** The GPS coordinates of the target delivery location and the estimated remaining distance.

This structured representation allows the drone to make informed decisions regarding navigation and obstacle avoidance [previous task results 1].

Action Space (A)

The action space available to the drone can be defined as follows:

- **Discrete Actions:** A set of predefined maneuvers, including:
 - Move forward
 - Ascend
 - Descend
 - Turn left
 - Turn right
- **Continuous Actions:** This space allows the drone to navigate freely in XYZ coordinates for more precise control. Continuous action control can be facilitated through algorithms such as Deep Deterministic Policy Gradient (DDPG), essential for sophisticated UAV navigation in complex urban settings [previous task results 2].

Reward Function (R)

The reward function needs to consider both the efficiency and safety of deliveries:

- **Positive Rewards:**
 - A significant reward for successfully reaching the delivery location (+100 points).
 - Incremental positive rewards for every step closer to the destination.
- **Negative Rewards:**
 - Penalties for collisions with obstacles (severity-based) and excessive battery usage.
 - A time penalty that deducts points for delays in reaching the delivery target.

Thus, the overall reward function can be expressed as:

$R = \text{Reward}$

This formulation promotes safe and efficient drone operation [previous task results 3]

2.

In reinforcement learning (RL), the exploration-exploitation dilemma is a critical aspect where an agent must decide between exploring new actions or exploiting known actions to maximize immediate rewards. Exploration refers to trying new actions to gain more information about the environment, while exploitation utilizes existing knowledge to secure the highest possible rewards from known actions.

The **ϵ -greedy strategy** is a common method to balance this trade-off. In this approach, the agent primarily chooses the action it believes has the highest value (exploitation) but with a small probability (ϵ), it randomly selects an action (exploration). This allows the agent to gather new information about less-explored parts of the action space without ignoring the potential of existing best actions.

For example, if $\epsilon = 0.1$, there is a 10% chance the agent will explore and a 90% chance it will exploit. This strategy effectively mitigates the risk of converging on suboptimal policies by ensuring that exploration occurs over time while facilitating the pursuit of optimal actions [previous task results 4, 10:230].

3.

To find the expected value of state s using the Bellman equation, we consider the information provided:

With a probability of 0.4, the agent transitions to state s' , yielding:

- Reward $R=10$
- Value $v(s')=5$

With a probability of 0.6, the agent transitions to state s'' , yielding:

- Reward $R=2$
- Value $v(s'')=3$

Applying the Bellman equation, the expected value of state s can be calculated as follows:

$$v(s) = \sum s' P(s'|s) (R + \gamma v(s'))$$

Substituting the values, where the discount factor $\gamma=0.5$:

$$v(s) = 0.4 \cdot (10 + 0.5 \cdot 5) + 0.6 \cdot (2 + 0.5 \cdot 3) = 0.4 \cdot (10 + 2.5) + 0.6 \cdot (2 + 1.5) = 0.4 \cdot 12.5 + 0.6 \cdot 3.5 = 5 + 2.1 = 7.1$$

Thus, the expected value of state s is 7.1. This computation highlights the agent's forecast of future rewards based on its current actions and state transitions.

In summary, defining the components of the MDP, discussing the exploration-exploitation trade-off, and calculating expected values using the Bellman equation are vital for developing an effective RL-based autonomous drone delivery system.