

Wids-2025 End-Term Report

FlappyBird: where Flappy learns to fly

Prashant Meena (24B0949)

31/01/2026

CHAPTER 1: WEEK 1&2&3&4

JN WEEK-1, I practiced Python and its libraries, including NumPy, Pygame, and Matplotlib.

In week-2, I studied Markov chains and Markov Decision Processes (MDP), focusing on state design, reward design, discretization, and environment testing.

The assignment solutions were uploaded to the GitHub repository: <https://github.com/prashantdlp/FlappyBird-AI>

MARKOV CHAINS AND MARKOV DECISION PROCESSES

Many sequential decision-making problems can be modeled as a Markov Decision Process (MDP) under the assumption of full state observability. An MDP is defined as a 5-tuple (S, A, P, R, γ) , where S is the set of states, A is the set of actions, P denotes the state transition probabilities (i.e., the consequences of actions), R is the reward function providing feedback to the agent, and $\gamma \in [0, 1]$ is the discount factor.

The discount factor γ determines the importance of future rewards. When $\gamma = 0$, the agent considers only immediate rewards, leading to greedy behavior. As γ approaches 1, future rewards are weighted more heavily, encouraging long-term planning. In episodic tasks, $\gamma = 1$ implies that future and immediate rewards are valued equally.

A BIT FORMAL :)

At each time step t , the agent observes the current state $s_t \in S$ and selects an action $a_t \in A$ according to a policy $\pi(a | s) = P(a_t = a | s_t = s)$. The environment then transitions to a new state s_{t+1} sampled from the transition distribution $P(\cdot | s_t, a_t)$ and emits a scalar reward $r_t = R(s_t, a_t)$.

The return is defined as the discounted cumulative reward:

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}.$$

The value function of a policy π is given by:

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

and the corresponding action-value function is:

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

The value function satisfies the Bellman expectation equation:

$$V^\pi(s) = \sum_{a \in A} \pi(a | s) \sum_{s' \in S} P(s' | s, a) [R(s, a) + \gamma V^\pi(s')].$$

The optimal value function is defined as:

$$V^*(s) = \max_\pi V^\pi(s),$$

and satisfies the Bellman optimality equation:

$$V^*(s) = \max_{a \in A} \sum_{s' \in S} P(s' | s, a) [R(s, a) + \gamma V^*(s')].$$

The objective of reinforcement learning is to find an optimal policy π^* that maximizes the expected discounted return:

$$\pi^* = \arg \max_\pi \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \right].$$

MARKOV PROPERTY

A state signal that succeeds in retaining all relevant information is said to be Markov, or to have the Markov property .

DISCRETIZATION

Tabular Q-learning requires a finite set of states, whereas most real-world environments are continuous. Discretization addresses this mismatch by mapping continuous variables into a finite number of bins. This introduces an inherent trade-off. Very fine discretization leads to a large Q-table, slow convergence, and sparse updates, while overly coarse discretization results in faster learning at the cost of precision and suboptimal decisions. The objective is to choose a resolution that is sufficient for effective decision-making without unnecessary complexity.

Flappy Bird Insight. In Flappy Bird, the agent does not require exact pixel-level information. Instead, decision-relevant abstractions are sufficient:

- Whether the bird is above or below the pipe gap
- The direction and magnitude of its vertical velocity
- Whether the next pipe is near or far

Thus, discretization should capture task-relevant structure rather than perfect environmental accuracy.

REWARD DESIGN

Reward design plays a crucial role in shaping agent behavior. The reward function should align closely with the true objective of the task while remaining simple and stable. Poorly designed rewards can lead to unintended behaviors, such as reward hacking or learning shortcuts that do not solve the intended problem. Sparse rewards make learning difficult due to limited feedback, whereas overly dense or large rewards can dominate learning and cause instability. Effective reward design balances guidance and flexibility, encouraging progress toward the goal without over-constraining the policy.

ENVIRONMENT VALIDATION AND STRESS TESTING

Before training an agent, it is essential to verify that the environment is well-defined and learnable. Random policies provide a simple yet effective sanity check, helping identify implementation bugs, impossible objectives, or poorly scaled rewards. If a random agent fails immediately, the task may be overly difficult; if it survives indefinitely, the task may be trivial. Unbounded or exploding rewards often indicate reward design errors.

Environment Stress Testing. Reinforcement learning systems can be highly sensitive to small changes in environment parameters. Stress testing evaluates robustness by varying key factors such as gravity, reward scaling, and observation noise. Agents that perform well only under a narrowly tuned configuration are brittle and lack generalization. Analyzing sensitivity to such perturbations provides insight into which environmental variables dominate the learned behavior.

VALUE FUNCTIONS

In reinforcement learning, the agent doesn't learn rules or strategies directly. Instead, it learns VALUE FUNCTIONS.

There are two important value functions :

1. State Value Function :- expected future reward when starting from state s
2. Action-value function :- expected future reward when starting from state s and taking action a

TEMPORAL DIFFERENCE (TD) LEARNING

Temporal Difference (TD) learning is a class of reinforcement learning methods that update value estimates incrementally, without waiting

for an episode to terminate. The central idea is to learn from the discrepancy between the predicted value of a state-action pair and a better estimate obtained after observing the next transition. This discrepancy is known as the *TD error*.

In Q-learning, the TD target is defined as

$$\text{target} = r + \gamma \max_{a'} Q(s', a'),$$

where r is the immediate reward obtained after taking action a in state s , s' is the resulting next state, and $\gamma \in [0, 1]$ is the discount factor that determines the importance of future rewards. The term $\max_{a'} Q(s', a')$ represents the estimated value of acting optimally from the next state onward.

Rather than directly replacing the current estimate, the Q-value is updated gradually using a learning rate $\alpha \in (0, 1]$:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(\text{target} - Q(s, a)).$$

This update moves the current estimate toward the TD target, allowing stable learning from partial experience and enabling online, step-by-step improvement.

Q-LEARNING

Q-learning is one of the most fundamental and widely used reinforcement learning algorithms. Its update rule is given by

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left(r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right),$$

where α is the learning rate controlling how strongly new information overrides old estimates, and γ is the discount factor that determines the importance of future rewards.

At first glance, this update rule looks familiar because it is directly derived from Temporal Difference (TD) learning. The term

$$r + \gamma \max_{a'} Q(s', a')$$

serves as the TD target, while the difference between this target and the current estimate $Q(s, a)$ represents the TD error.

Q-learning is particularly powerful because it is *model-free*, *off-policy*, and *tabular*. It does not require any knowledge of the environment's transition dynamics or reward probabilities. Instead, it learns purely from experience in the form of observed tuples (s, a, r, s') . As long as an environment can provide states, actions, rewards, and next states, the same Q-learning algorithm can be applied without modification.

As a result, Q-learning can be used across a wide variety of environments such as

multi-armed bandits, GridWorld, FrozenLake, Taxi-v3, and even games like Flappy Bird. The environment changes, but the learning algorithm remains exactly the same—highlighting the generality and elegance of Q-learning.

EXPLORATION VS. EXPLOITATION

A reinforcement learning agent often faces a fundamental dilemma: should it choose the best-known action based on current knowledge, or should it try a new action to discover potentially better outcomes? This trade-off is known as the *exploration vs. exploitation* problem. Pure exploitation can cause the agent to converge to a sub-optimal policy if better actions are never tried, while pure exploration prevents convergence altogether, as the agent never commits to the best-performing actions.

ϵ -GREEDY STRATEGY

A simple and effective way to balance exploration and exploitation is the ϵ -greedy strategy. With probability ϵ , the agent selects an action uniformly at random (exploration), and with probability $1 - \epsilon$, it selects the action with the highest estimated Q -value (exploitation). To encourage efficient learning, ϵ is typically decayed over time, allowing the agent to explore more during early training and gradually shift toward exploitation as its value estimates become more accurate.

PSEUDO CODE EXAMPLE

Algorithm 1 Q-learning for Flappy Bird

```
1: Initialize Q-table  $Q(s, a)$  arbitrarily
2: Set learning rate  $\alpha$ , discount factor  $\gamma$ 
3: Set exploration rate  $\epsilon$  and decay factor  $d$ 
4: for each episode do
5:   Reset environment and observe initial
     state  $s$ 
6:   while episode not terminated do
7:     With probability  $\epsilon$ , select a random
       action  $a$ 
8:     Otherwise, select  $a = \arg \max_{a'} Q(s, a')$ 
9:     Execute action  $a$ 
10:    Observe reward  $r$  and next state  $s'$ 
11:    Update  $Q$ -value:
```

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

```
12:   Set  $s \leftarrow s'$ 
13: end while
14: Decay exploration rate:  $\epsilon \leftarrow d \cdot \epsilon$ 
15: end for
```