

---

# **Report on Kalman Filter applied on Neural Network for Finance**

---

August 13, 2025

Mathys VINATIER - [GitHub Project page](#)

Supervisor:

Pr Kim Tae-Wan

Mathys VINATIER

# Contents

10 Week 10	1
------------	---

# Chapter 10

## Week 10

### Contents

---

<b>10.1 Reminder on Greedy Policy and QLearning . . . . .</b>	<b>2</b>
<b>10.2 Implementation Progress and Experiments . . . . .</b>	<b>3</b>
10.2.1 Environment Setup . . . . .	3
10.2.2 Preliminary Results for Q-Learning Greedy Policy . . . . .	5
<b>10.3 Challenges and Solutions . . . . .</b>	<b>14</b>
<b>10.4 Next Steps . . . . .</b>	<b>14</b>

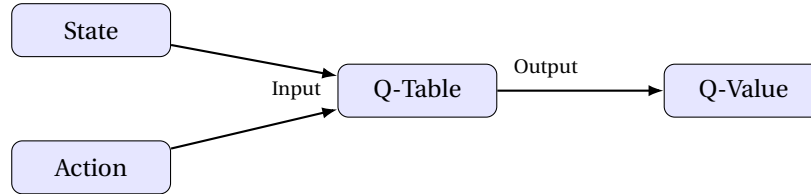
---



## 10.1 Reminder on Greedy Policy and QLearning

Q-Learning is a **model-free, off-policy, value-based reinforcement learning algorithm** that trains an action-value function (Q-function) to find the optimal policy indirectly. The Q-function estimates the expected cumulative reward of taking a certain action in a given state and following the best policy thereafter.

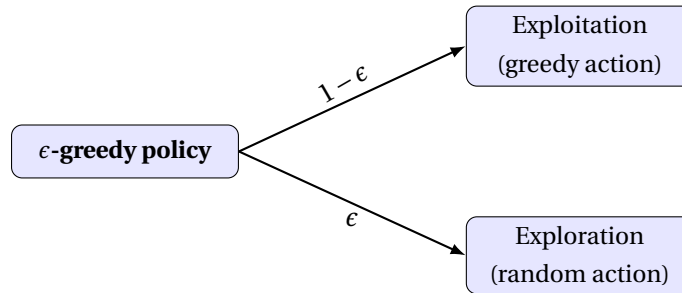
- **Q-table:** Internally, Q-Learning uses a Q-table storing values for each state-action pair. Initially, all values are zero, and the table is updated iteratively during training.



- **Epsilon-Greedy Strategy:** To balance exploration and exploitation, actions are chosen using an epsilon-greedy policy:

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

- With probability  $1 - \epsilon$ , exploit by selecting the action with the highest Q-value.
- With probability  $\epsilon$ , explore by selecting a random action.
- $\epsilon$  decays over time to shift from exploration to exploitation.



- **TD Update:** At each step, the Q-value for the current state-action pair  $Q(s_t, a_t)$  is updated based on the immediate reward  $r_{t+1}$  plus the discounted maximum Q-value of the next state  $s_{t+1}$ :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

- **Off-policy Learning:** The policy used to select actions (epsilon-greedy) differs from the policy used to compute the target update (greedy). This distinction allows Q-Learning to learn optimal policies even when exploring randomly.

After enough training episodes, the Q-table converges to the optimal Q-function, which directly yields the optimal policy by choosing actions with the highest Q-values in each state.

## 10.2 Implementation Progress and Experiments

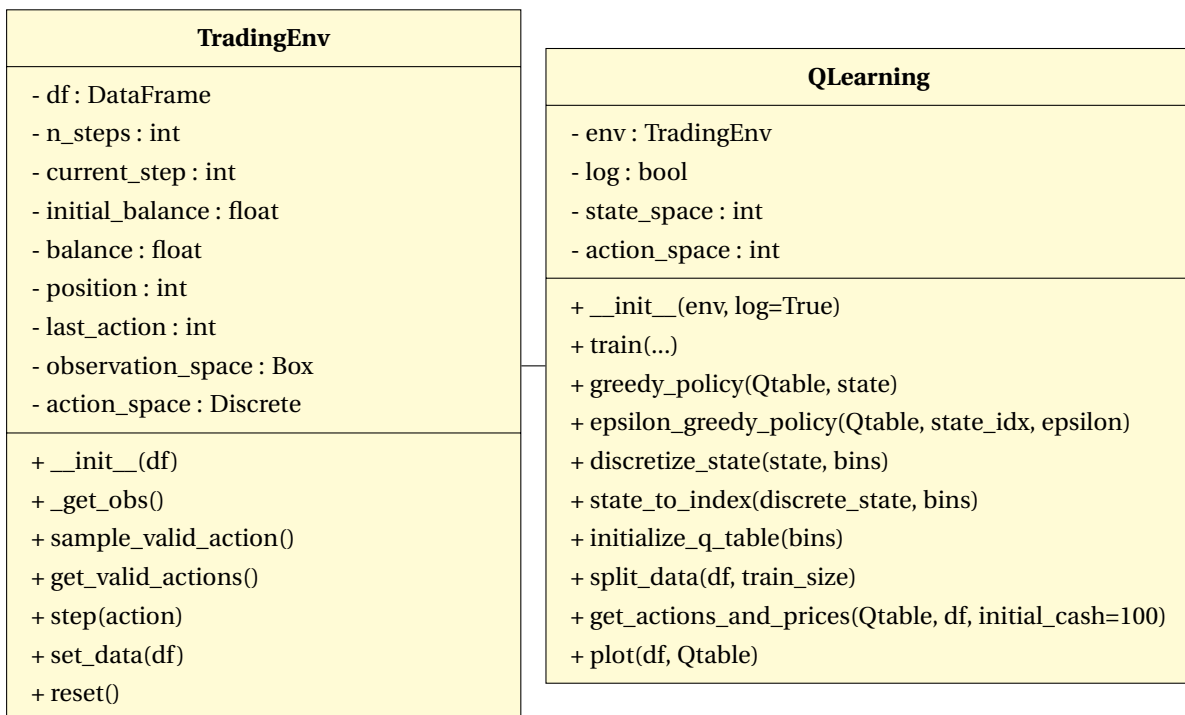
This week was dedicated to the initial implementation and experimentation phase of the project, focusing primarily on the Proximal Policy Optimization (PPO) basis algorithms. Using the Hugging Face classes, we developed a prototype with greedy policy agent and integrated it into a custom trading environment inspired by the concepts discussed in previous weeks.

### 10.2.1 Environment Setup

The custom environment was designed to simulate stock market trading dynamics with the following key elements :

- **State representation** including historical price data, technical indicators (moving averages, RSI) and Kalman filtered signals.
- **Action space** consisting of discrete trading actions: Buy, Hold or Sell.
- **Reward structure** based on portfolio value change and risk-adjusted metrics to incentivize both profit and stability.
- **Episode design** with fixed-length trading periods to allow episodic evaluation of agent performance.

The environment was implemented adhering to OpenAI Gym interfaces, allowing smooth integration with standard RL training pipelines. Here is the UML of the created environnement :



**Algorithm 1: Trading Environment Behavior****Input:** Environment  $env$ , number of episodes  $N$ , exploration rate  $\epsilon$ **Output:** Episode value and rewards

```

1 for  $episode \leftarrow 1$  to  $N$ 
2    $s \leftarrow env.reset()$ ;
3    $total\_reward \leftarrow 0$ ;
4   while  $True$ 
5     Get valid actions  $A_{valid} \leftarrow env.get\_valid\_actions()$ ;
6     Select  $a \leftarrow$  random choice from  $A_{valid}$ ;
7     Execute  $(s', r, done, info) \leftarrow env.step(a)$ ;
8      $total\_reward \leftarrow total\_reward + r$ ;
9     Log( $episode, s, a, r, info["portfolio\_value"]$ );
10     $s \leftarrow s'$ ;
11    if  $done$  then
12      break
13  Print episode summary :  $total\_reward$ , final portfolio value;

```

**Algorithm 2: Q-Learning Training****Input:** Environment  $env$ , data frame  $df$ , training size  $train\_size$ , episode  $N$ **Output:** Trained Q-table  $Q$ 

```

1 Split data:  $df_{train} \leftarrow train\_size$  of  $df$ ;
2 Calculate bins for discretization :  $bins \leftarrow$  compute bins from  $df_{train}$ ;
3 Initialize Q-table
4 for  $episode \leftarrow 1$  to  $N$ 
5   Exploration rate :  $\epsilon \leftarrow \epsilon_{min} + (\epsilon_{max} - \epsilon_{min}) \cdot \exp(-decay\_rate \times episode)$ ;
6   Reset environment :  $state_{cont} \leftarrow env.reset()$ ;
7   Discretize initial state :  $state_{disc} \leftarrow discretize(state_{cont}, bins)$ ;
8   Convert to index :  $state_{idx} \leftarrow state\_to\_index(state_{disc}, bins)$ ;
9   for  $step \leftarrow 1$  to  $max\_steps$ 
10    Choose action  $a$  using epsilon-greedy policy :

$$a \leftarrow \begin{cases} \arg\max_{a'} Q[state_{idx}, a'] & \text{with probability } 1 - \epsilon \\ \text{random valid action} & \text{with probability } \epsilon \end{cases}$$

11    Take action :  $(next_{state_{cont}}, r, done, info) \leftarrow env.step(a)$ ;
12    Discretize next state :  $next_{state_{disc}} \leftarrow discretize(next_{state_{cont}}, bins)$ ;
13    Convert to index :  $next_{state_{idx}} \leftarrow state\_to\_index(next_{state_{disc}}, bins)$ ;
14    Update Q-value
    :  $Q[state_{idx}, a] \leftarrow Q[state_{idx}, a] + \alpha (r + \gamma \max_{a'} Q[next_{state_{idx}}, a'] - Q[state_{idx}, a])$ ;
15     $state_{idx} \leftarrow next_{state_{idx}}$ ;
16    if  $done$  then
17      break;
18 return  $Q$ 

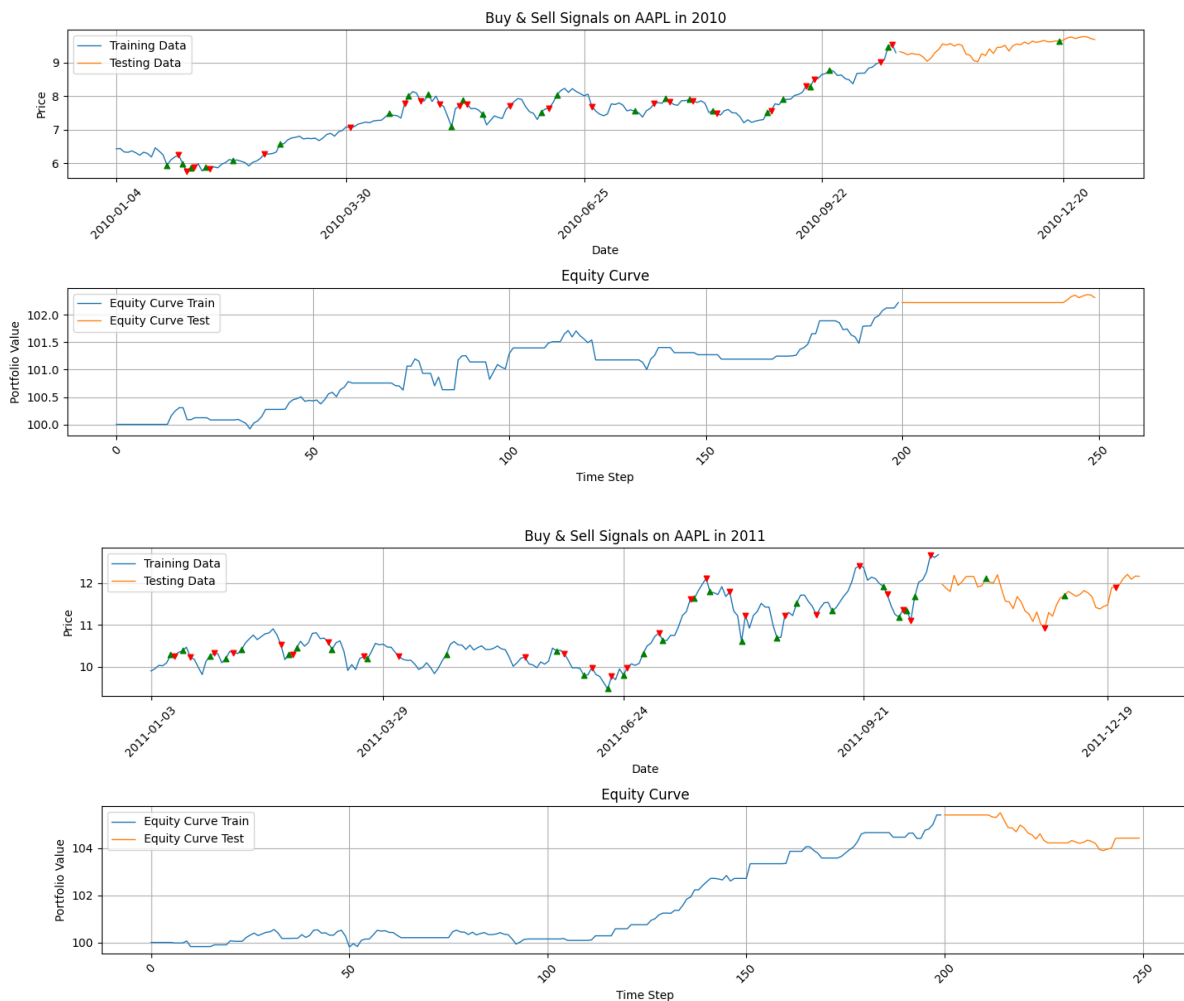
```

In the context of Q-learning for trading environments, bins are used to discretize the continuous state space into a finite number of discrete states. This discretization is essential because the Q-table requires a finite and manageable number of states to store and update the expected rewards for each state-action pair. Continuous variables, such as stock prices or technical indicators, have infinite possible values, which makes it impractical to represent them directly in a Q-table. To address this, we divide the range of each continuous feature into a fixed number of intervals called bins. Each bin represents a segment of the feature's value range, allowing us to map any continuous observation into a discrete category. For implementation, we analyze the training dataset and create evenly spaced bins (10 bins) for each feature by computing linearly spaced intervals between the minimum and maximum values of that feature. Then, during training or evaluation, the observed continuous state is converted into a discrete state by determining which bin each feature value falls into. This discretized state is subsequently converted into a unique index to access and update the Q-table. The use of bins thus enables effective Q-learning in continuous state environments by simplifying the state representation while preserving essential information.

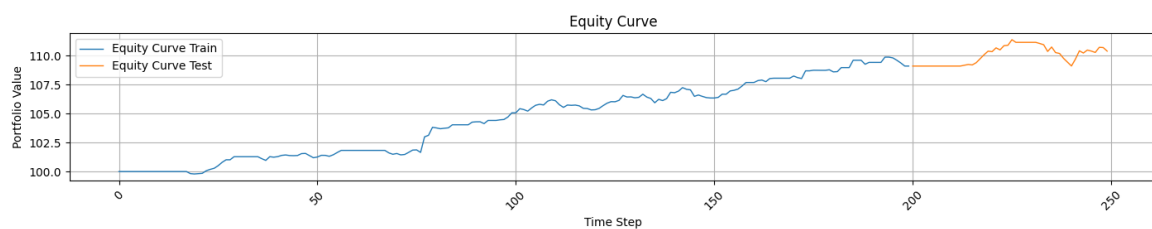
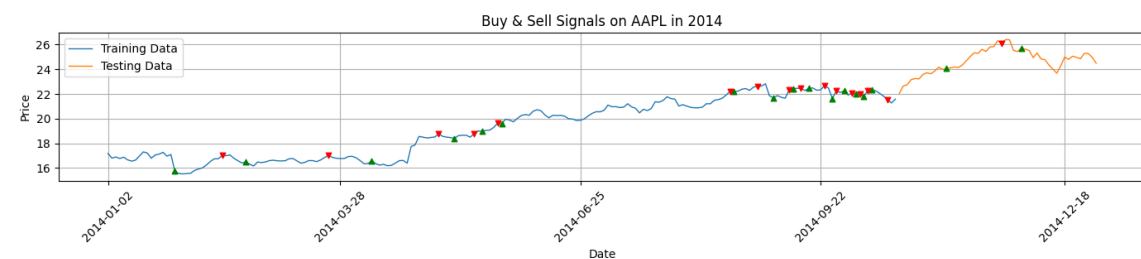
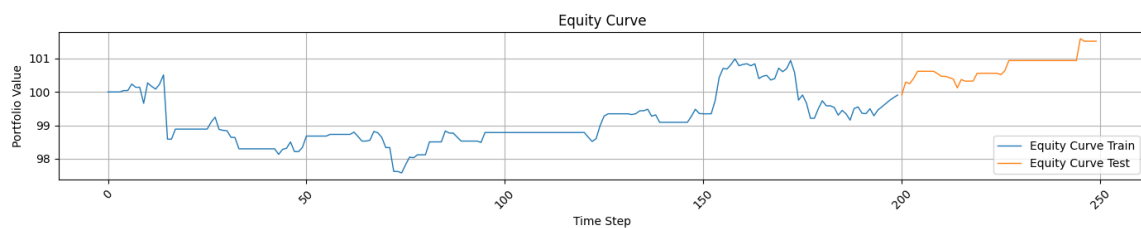
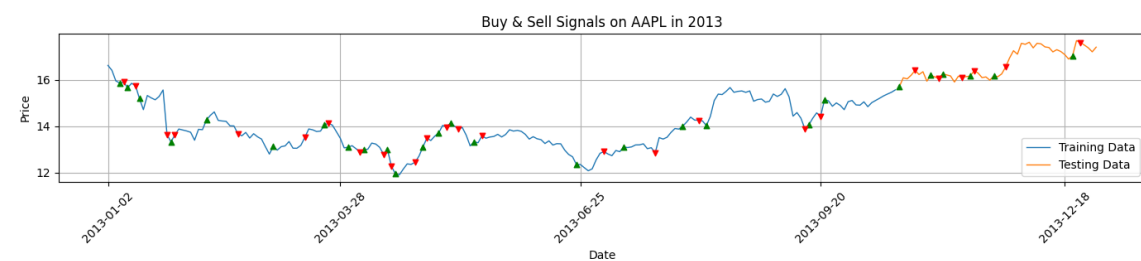
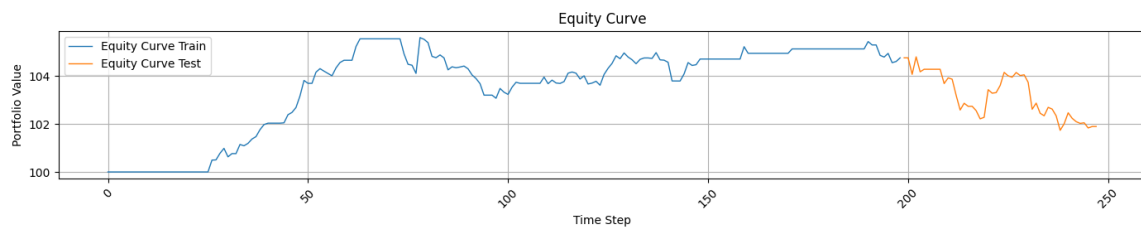
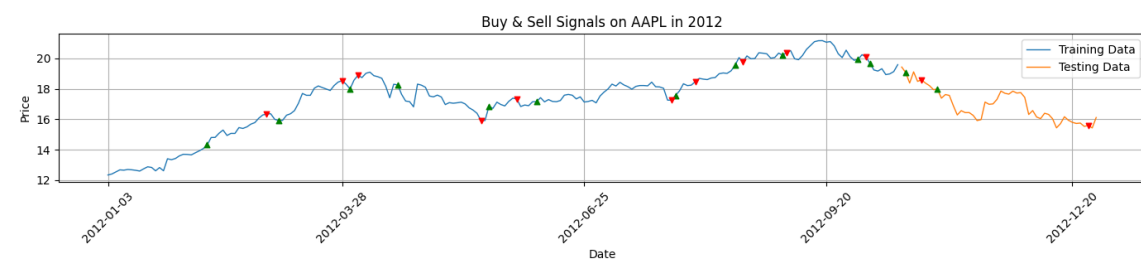
### 10.2.2 Preliminary Results for Q-Learning Greedy Policy

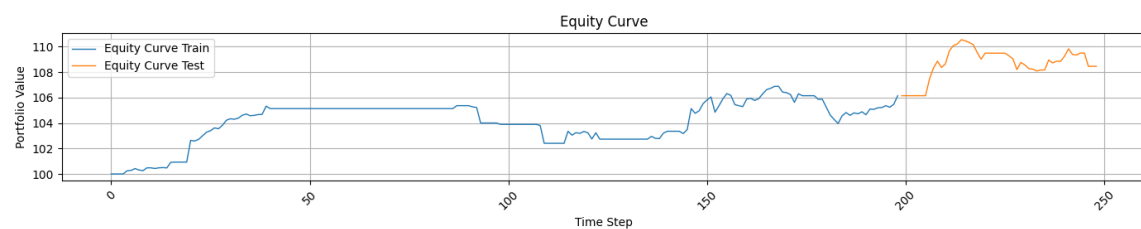
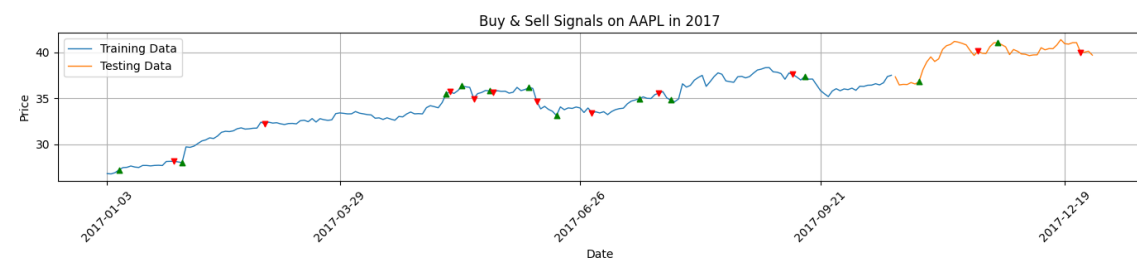
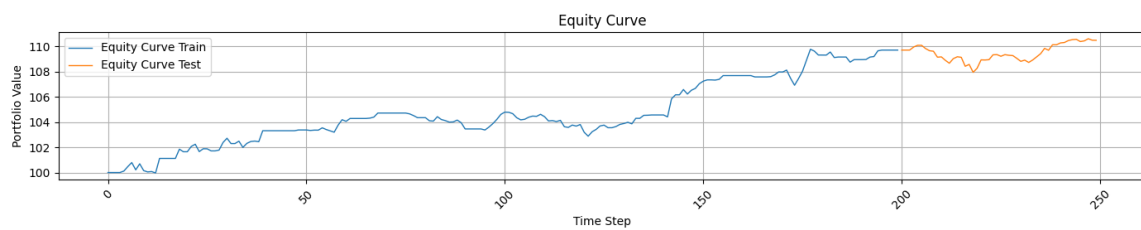
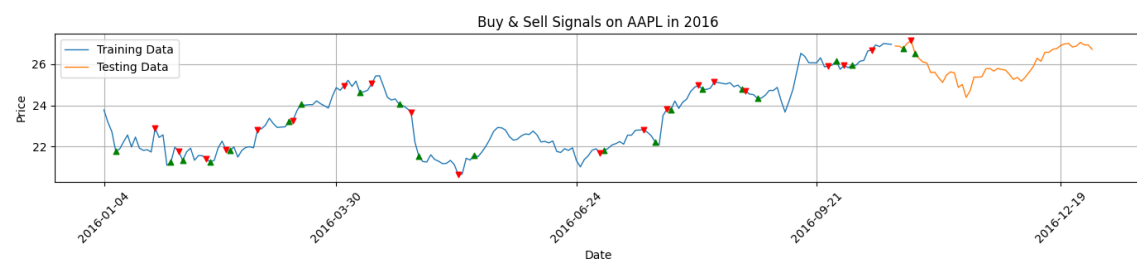
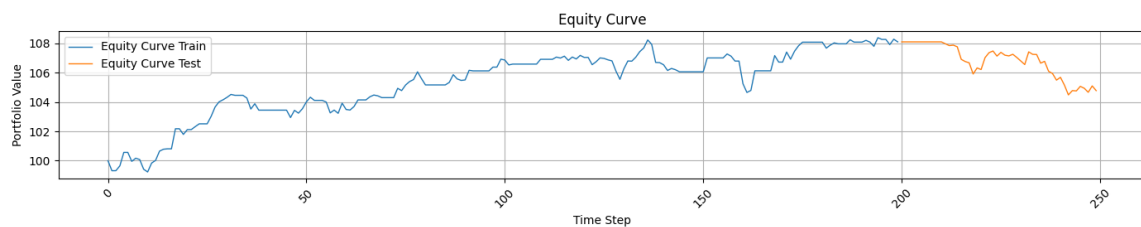
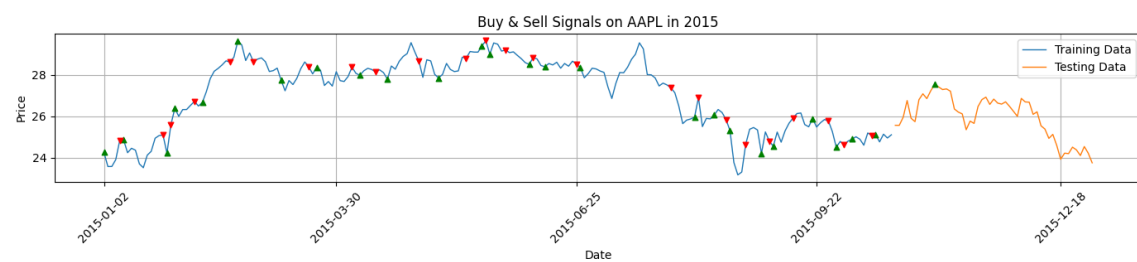
Initial experiments aimed to evaluate the training behavior and reward progression of the Q-Learning agent following a greedy policy. Results were compared over two time horizons : a shorter training period of 100 episodes and an extended training of 1000 episodes.

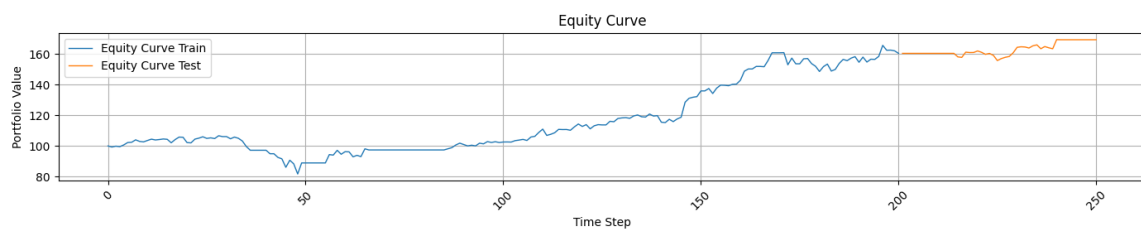
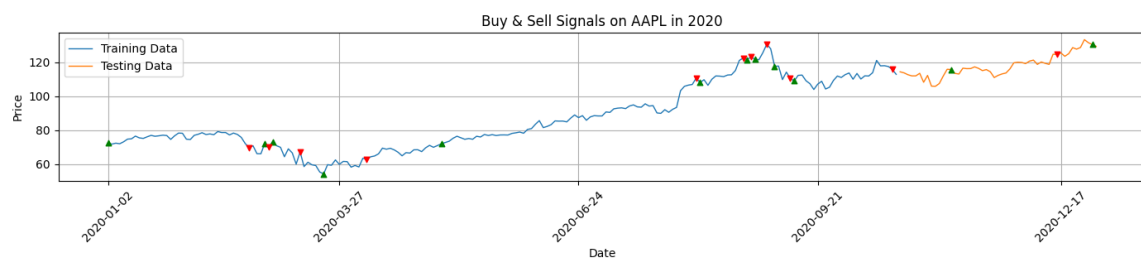
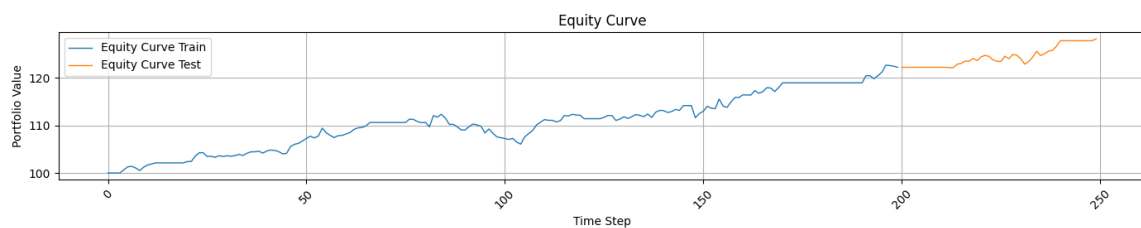
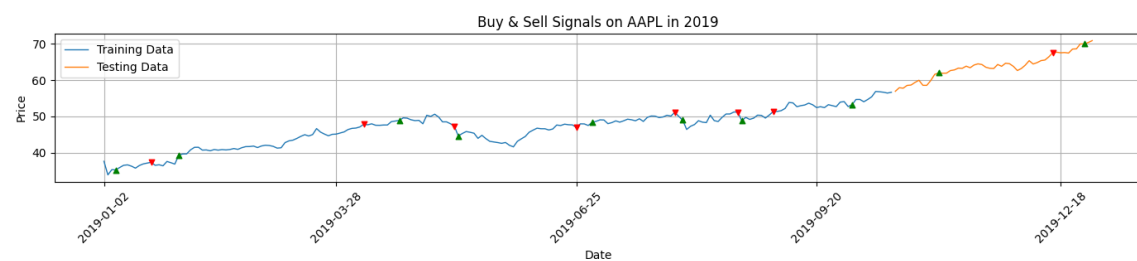
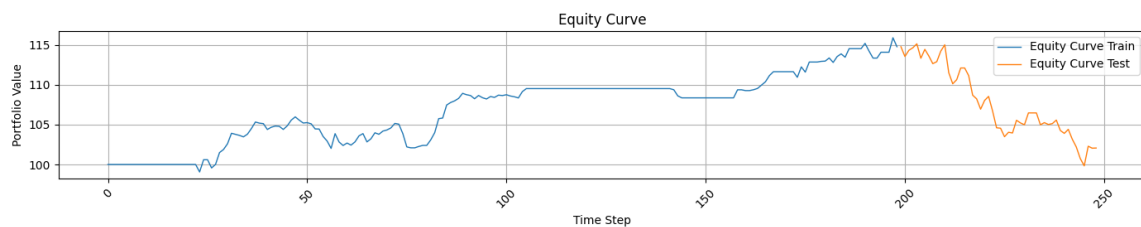
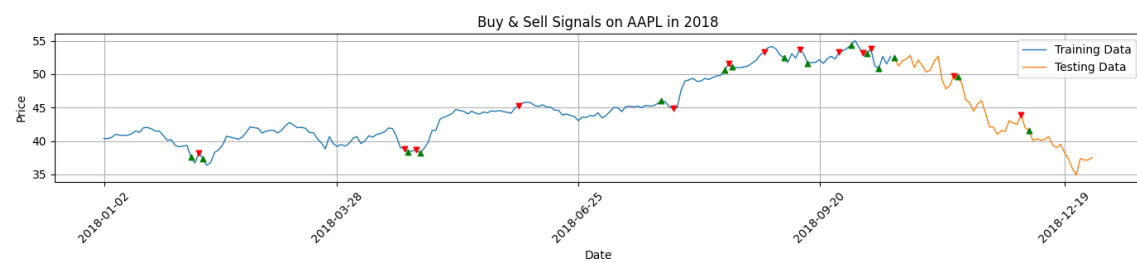
#### 100 episodes

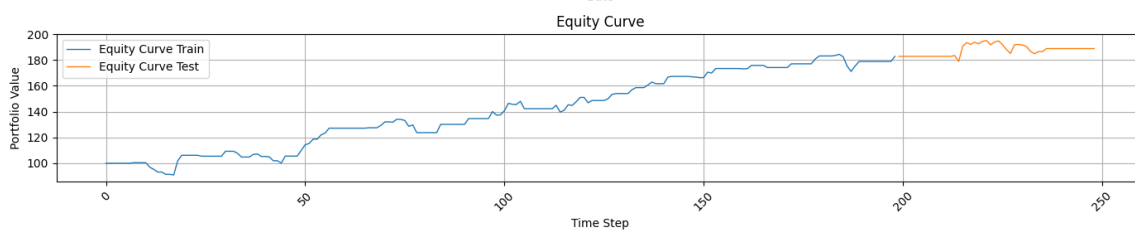
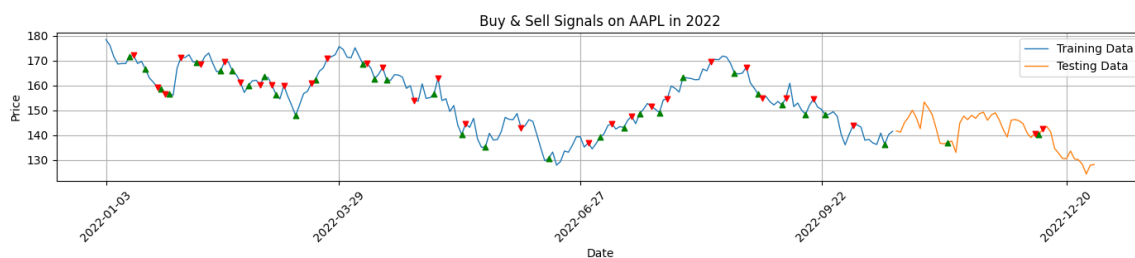
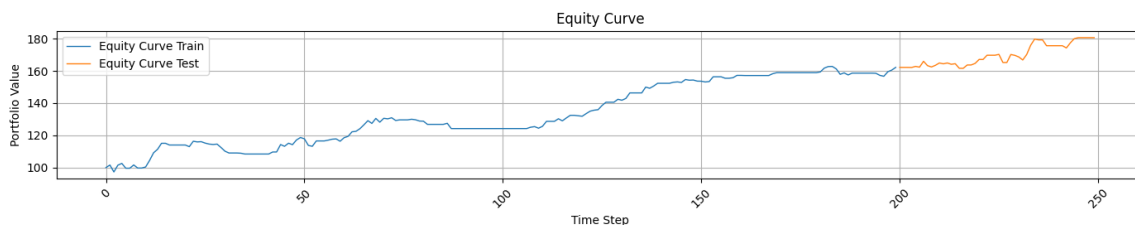
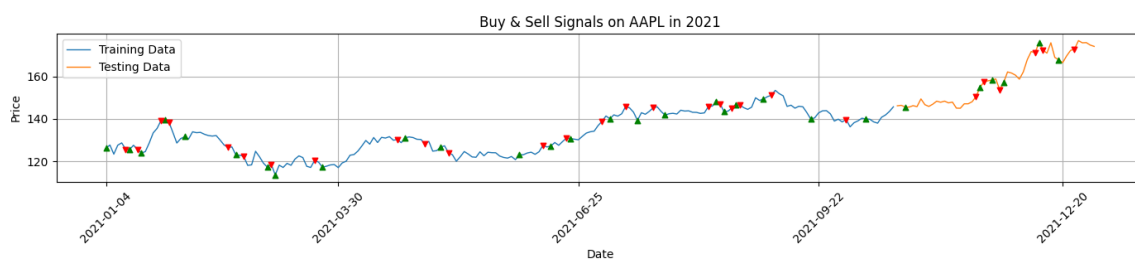




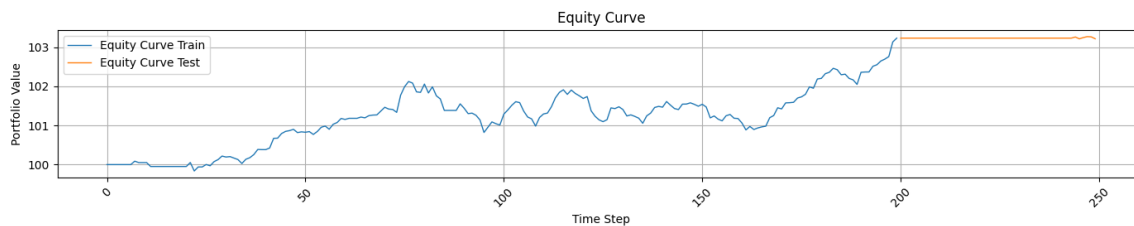
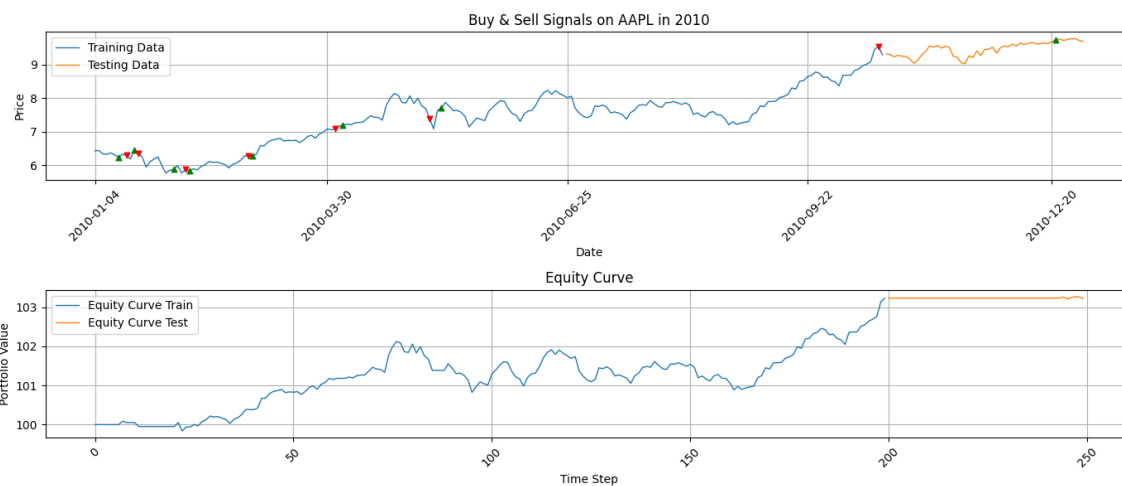


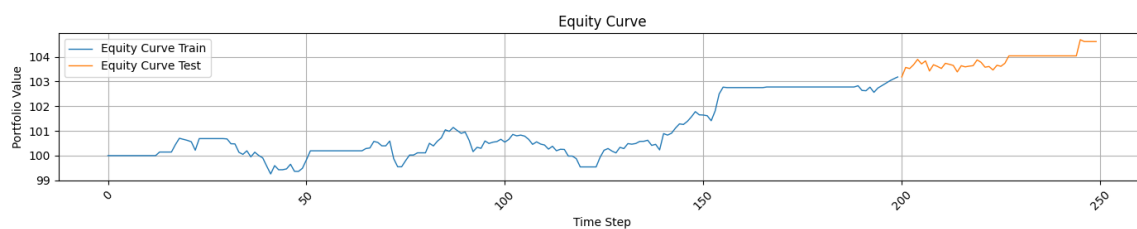
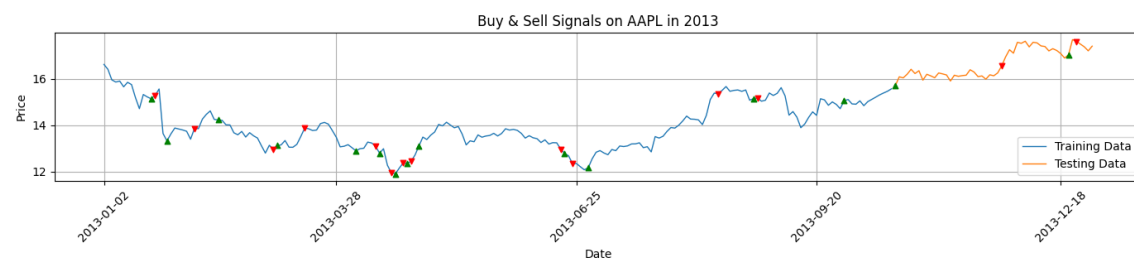
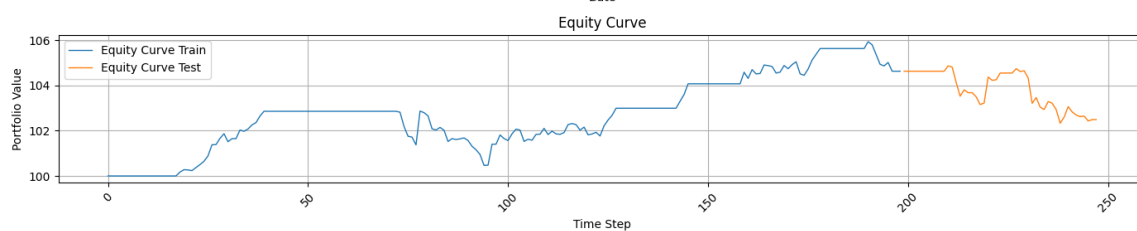
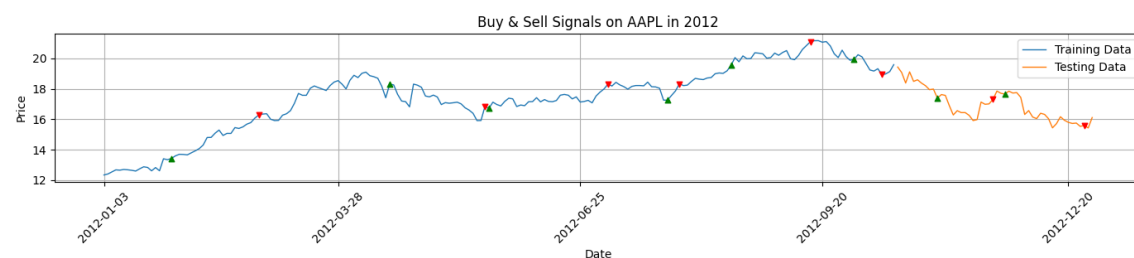
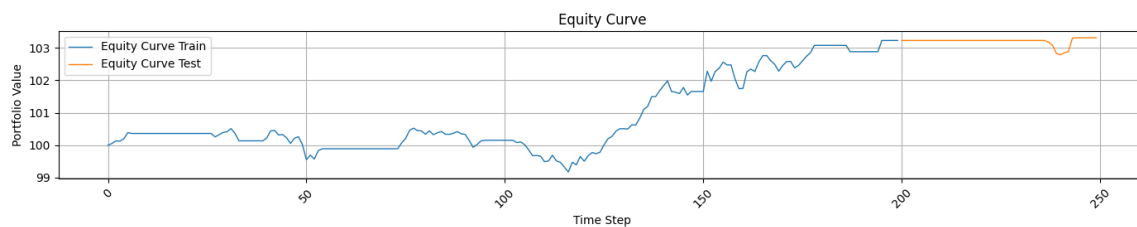
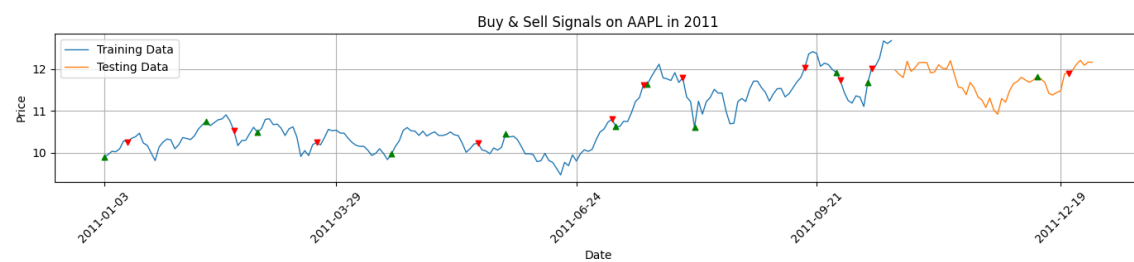


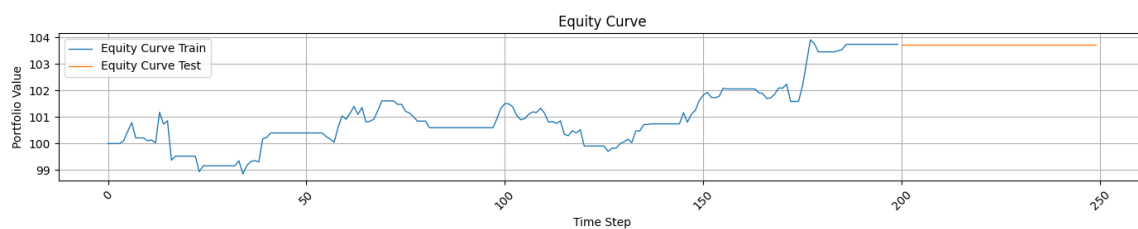
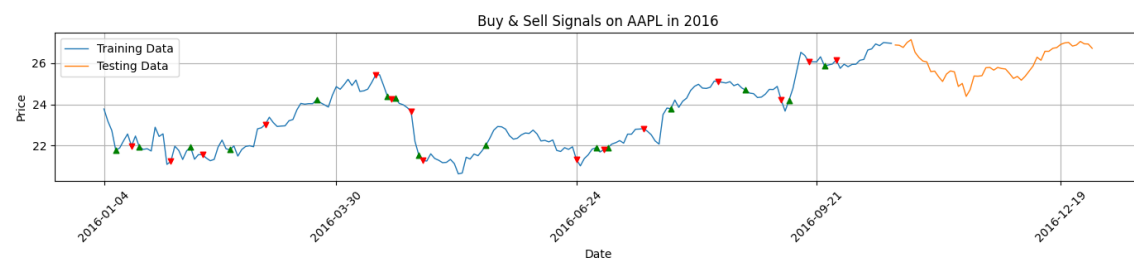
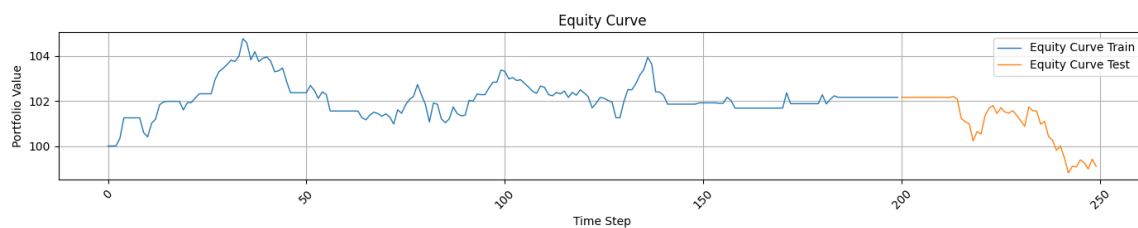
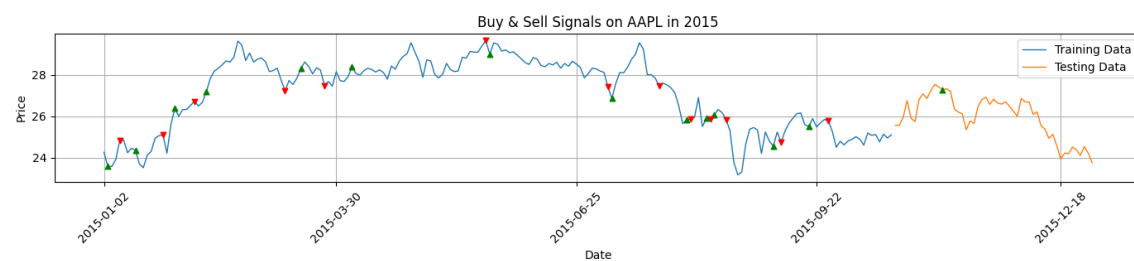
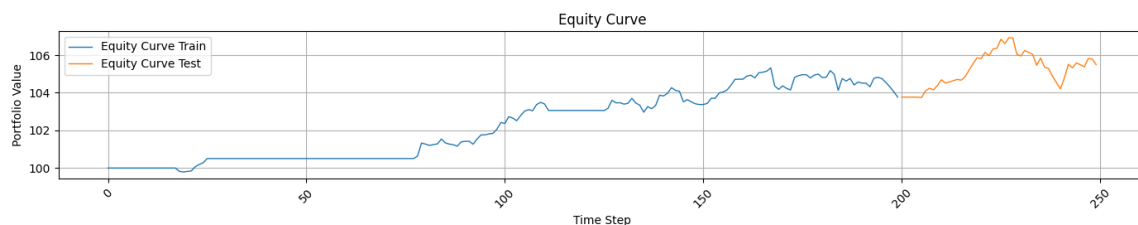
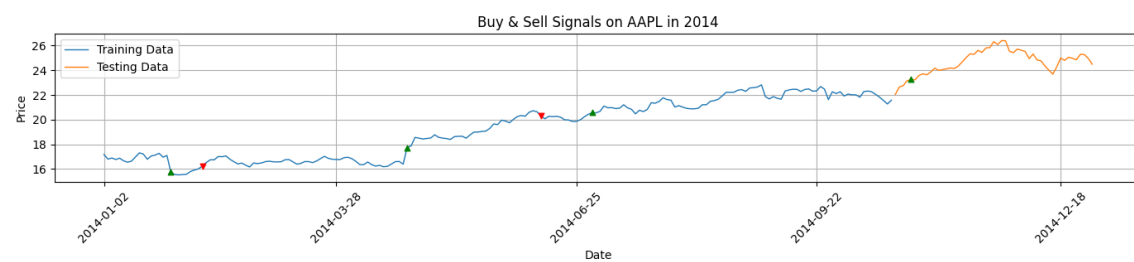


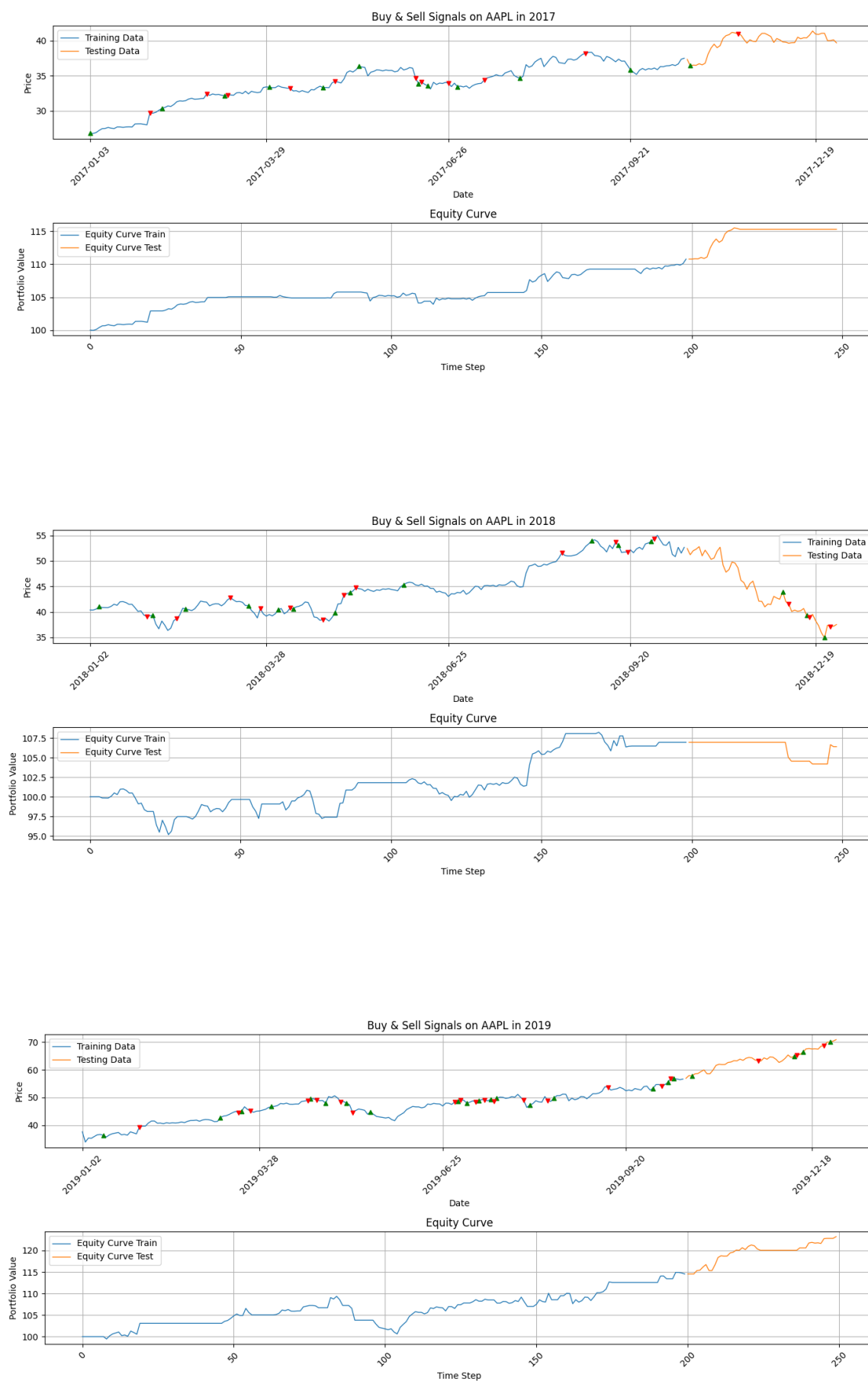


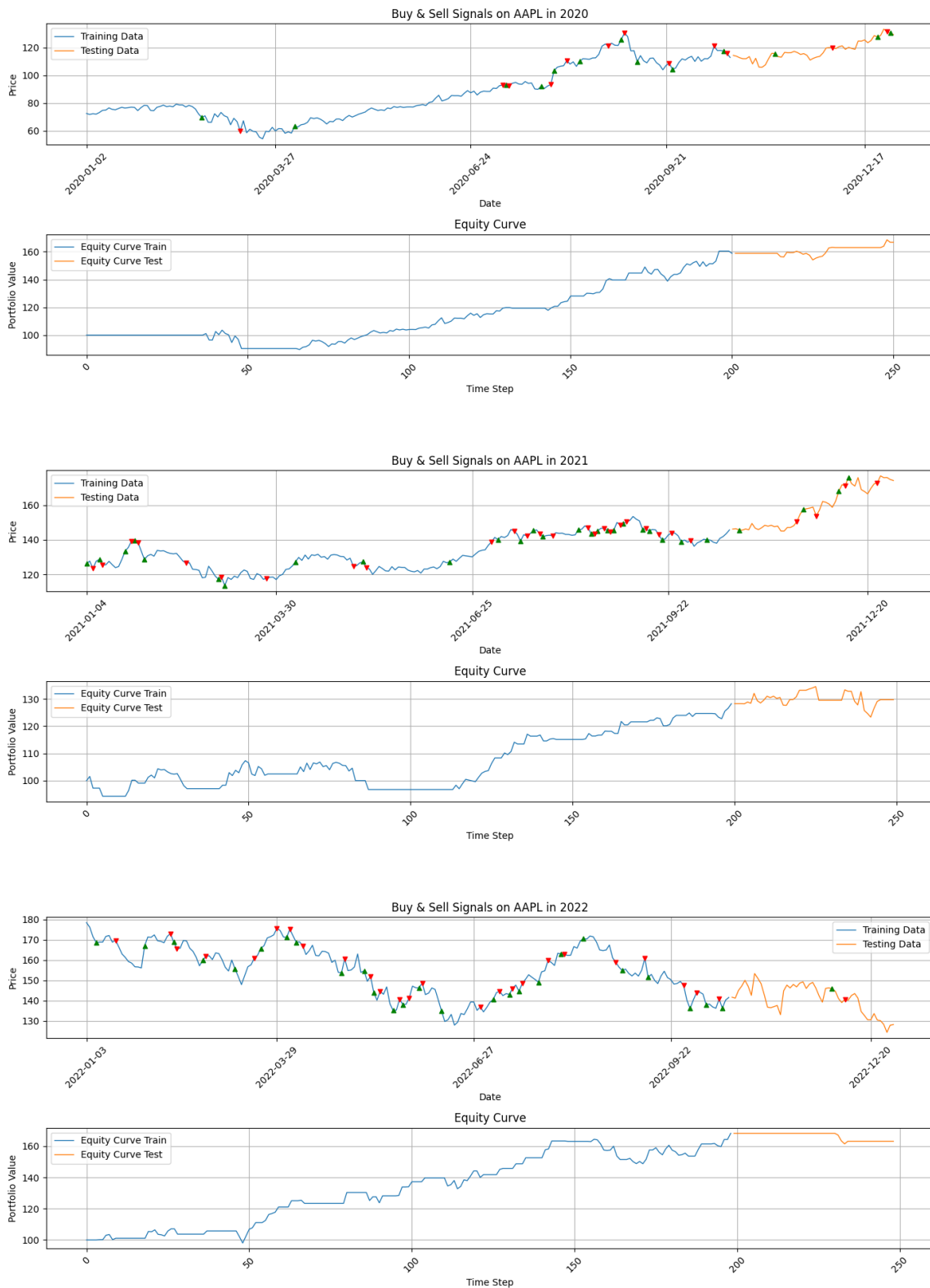
## 1000 episodes











- Over 100 episodes, the cumulative rewards exhibited a modest upward trend, indicating initial learning progress, though with significant variability.
- Extending to 1000 episodes resulted showed that the agent is not taking position anymore, a convergence of the hold position should be evaluated
- Value estimates became more consistent over time, though some fluctuations persisted due to the noisy environment dynamics (should add filtering or smoothing).



- The results highlight the trade-off between exploration and exploitation inherent in the greedy approach and motivate future exploration of strategies with controlled exploration. Noticed that the agent is not learning from the path but more on the direct previous value.

These preliminary findings establish a baseline for Q-Learning performance and set the stage for comparisons with other reinforcement learning methods such as Deep Q-Networks (DQN).

### 10.3 Challenges and Solutions

The Q-Learning experiments faced several challenges specific to the greedy policy and the trading environment :

- **Limited exploration due to greediness** : The strictly greedy policy limited the agent's ability to discover better actions, particularly early in training. Future work will include incorporating greedy or softmax action selection to balance exploration.
- **Reward noise and sparsity** : The stochastic environment produced noisy and sparse rewards, complicating learning.
- **Hyperparameter sensitivity** : Learning rate and discount factor tuning were critical, especially over longer training horizons. We will use optuna to get better hyperparameter in the future on more complexe models.
- **Computational efficiency** : Although less demanding than deep RL methods, Q-Learning with large state spaces required careful management of updates. We leveraged experience replay buffers and batch updates to accelerate convergence.

### 10.4 Next Steps

Building on the current progress with Q-Learning and the greedy policy, the immediate next steps focus on advancing the model complexity and preparing for PPO :

- Designing and integrating a deep neural network architecture to implement Deep Q-Learning (DQN), enabling function approximation for larger and continuous state spaces.
- Conducting experiments to compare the performance of DQN against the baseline Q-Learning greedy policy.
- Following the DQN implementation, proceeding with the development and fine-tuning of the PPO agent, including enhanced reward shaping and hyperparameter optimization.
- Expanding the trading environment to model more realistic market conditions such as transaction costs, slippage, and possibly varying liquidity.
- Performing ablation studies to isolate the effects of PPO-specific components like clipping, value function loss weighting, and entropy regularization.
- Enhancing visualization tools to track detailed performance metrics such as Sharpe ratio, maximum drawdown, and other risk-adjusted return measures.

These steps will establish a robust foundation for evaluating advanced reinforcement learning algorithms in the trading domain and inform subsequent research and publication efforts.

