



# **Trading with Deep Reinforcement Learning**

## **Description and Experiments with Synthetic Data**

Denis M. Becker

## Table of Contents

Table of Contents.....	2
1 Stable Baselines 3 and Open AI Gym.....	3
2 Environments .....	4
3 Experiments with Synthetic Time Series .....	10
4 Take-aways, To-Do's and Future research .....	49

## 1 Stable Baselines 3 and Open AI Gym

- We apply deep reinforcement **algorithms** provided by **Stable Baselines 3** (Python library that provides a set of pre-implemented RL algorithms)
- Particularly:

Continuous decision space	Discrete decision space
---	Deep Q-network (DQN)
Proximal Policy Optimization (PPO)	Proximal Policy Optimization (PPO)
Advantage Actor Critic (A2C)	Advantage Actor Critic (A2C)
Deep Deterministic Policy Gradient (DDPG)	---
Twin Delayed DDPG (TD3)	---
Soft Actor Critic (SAC)	---

## 2 Environments

- We consider **two** different **RL environments**:
  - **dichotomic** (binary) **decision space** «In» (buy or hold) and «out» (sell or stay-out)
  - **Continues actions** in interval [0,1]
- Both implemented by means of **Open AI Gym** (container that contains many useful helper functions and terminology used in SB3)
- The **RL agent maximizes**:
  - either total **wealth** in the end of one trading epoch,
  - or **risk-adjusted performance**.

### 2.1 Dichotomic Trading Actions

**Initialization:**

$B_{\text{Out}} = B_0, S_{\text{Out}} = 0, t = 0$ , Observe state  $S_0$ .

**Loop through all periods:**

The outgoing values from the previous iteration become the ingoing values:

$$B_{\text{In}} = B_{\text{Out}}$$

$$S_{\text{In}} = S_{\text{Out}}$$

Obtain new state  $S_t$  and receive the asset's return  $r_{\text{Asset},t}$

Dependent on  $S_t$  take the binary decision  $a \in \{0,1\}$

## 2 Environments

If  $a = 1$  (we go or stay long in stocks):

If  $S_{In} = 0$  (out of risky asset, and long in bonds, i.e.,  $B_{In} \geq 0$ ):

We need to buy risky asset: \*

$$S_{Out} = \frac{B_{In}}{(1 + \tau)} \cdot (1 + r_{Asset,t})$$

$$B_{Out} = 0$$

Else (i.e.  $S_{In} > 0$  and  $B_{In} = 0$  which means long in stocks, out of bonds):

We will stay in risky asset:

$$S_{Out} = S_{In} \cdot (1 + r_{Asset,t})$$

$$B_{Out} = 0$$

## 2 Environments

If  $a = 0$  (we sell stocks or stay out):

If  $S_{In} = 0$  (out of stocks, and long in bonds, i.e.,  $B_{In} \geq 0$ ):

$$S_{Out} = 0$$

$$B_{Out} = B_{In} \cdot (1 + r_B)$$

Else (i.e.,  $S_{In} > 0$  and  $B_{In} = 0$  which means long in stocks, out of bonds):

We need to sell stocks: \*

$$B_{Out} = S_{In} \cdot (1 - \tau) \cdot (1 + r_B)$$

$$S_{Out} = 0$$

### 2.2 Continuous Trading Actions

- Proportion/share of risky asset after trading and after transaction costs is denoted as  $x$ .
- Short sales are not allowed, i. e.  $x$  is continuous in the interval  $x = [0,1]$ .
- Transaction costs are a percentage of value of sold/bought stocks.

**Initialization:**

$B_{\text{Out}} = B_0$ ,  $S_{\text{Out}} = 0$ ,  $t = 0$ , Observe state  $S_0$ .

### Loop through all periods:

The outgoing values from the previous iteration become the ingoing values:

$$B_{\text{In}} = B_{\text{Out}}, S_{\text{In}} = S_{\text{Out}}$$

Obtain new state  $S_t$  and receive the asset's return  $r_{\text{Asset},t}$ .

Dependent on  $S_t$  take the continuous decision  $a \in [0,1]$

Transform the share in the risky asset, after transaction costs.

$$\Delta S = \frac{B_{\text{In}} - \left(\frac{1}{x} - 1\right) \cdot S_{\text{In}}}{\frac{1}{x} + \tau}$$

Calculate the outgoing values of the risky asset and bank account:

$$B_{\text{Out}} = (B_{\text{In}} - \Delta S \cdot (1 + \tau)) \cdot (1 + r_B)$$

$$S_{\text{Out}} = (S_{\text{In}} + \Delta S) \cdot (1 + r_S)$$

## 3 Experiments with Synthetic Time Series

- All experiments are run on a virtual Windows machine created in OpenStack with the following specifications:

Operating system:	Windows Server 2022 Standard, version 21H2, OS build: 20348.887
Processor:	Intel Xeon E312xx (Sandy Bridge, IBRS update) 2.60 GHz (2 processors)
Installed RAM:	32.0 GB
System type:	64-bit operating system, x64-based processor

### 3 Experiments

#### 3.1 Experiment #1 - Deterministic, repetitive sequence

- Time series has the following form:

Time	1	2	3	4	5	6	7	8	9	10	11	12	Rep
Return	2 %	-2 %	1 %	2 %	1 %	4 %	-1 %	5 %	-2 %	4 %	-1 %	-2 %	2 % ...

- Repeated four times when training:  $12 \times 4 = 48$  timesteps
- Optimal overall return when starting with  $B_0 = 10 \Rightarrow 8.34$
- Observation in  $t$  consists of only  $r_t$  (no information on earned balance, lagged returns, etc.)
- Environment #1: Continues action space in interval  $[0,1]$
- Environment #2: Dichotomic action space: 0 or 1

### 3 Experiments

- Transition probabilities, expected state returns and optimal actions:

	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	$\mathbb{E}[r_{t+1}]$	Optimal action
<i>from</i>	-2 %	-1 %	1 %	2 %	4 %	5 %		
<i>from</i>	-2 %			1/3	1/3	1/3	>0	Long
<i>from</i>	-1 %	1/2					>0	Long
<i>from</i>	1 %				1/2	1/2	>0	Long
<i>from</i>	2 %	1/2		1/2			<0	Out
<i>from</i>	4 %		1				<0	Out
<i>from</i>	5 %	1					<0	Out

### 3 Experiments

- Perfect information is given if states contain information  $(r_{t-2}, r_{t-1})$ . It can be represented by the following transition matrix.

	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>	<i>to</i>
<i>from</i>	2%, -2%	-2 %,	1 %	2 %	1 %	4 %	-1 %	5 %	-2 %	4 %	-1 %	-2 %	-2 %
	-2%	1 %	2 %	1 %	4 %	-1 %	5 %	-2 %	4 %	-1 %	-2 %	2 %	
<i>from</i>	2%, -2%		1										
<i>from</i>	-2 %, 1%			1									
<i>from</i>	1 %, 2%				1								
<i>from</i>	2 %, 1 %					1							
<i>from</i>	1 %, 4 %						1						
<i>from</i>	4 %, -1 %							1					
<i>from</i>	-1 %, 5%								1				
<i>from</i>	5 %, -2 %									1			
<i>from</i>	-2 %, 4 %										1		
<i>from</i>	4 %, -1 %											1	
<i>from</i>	-1 %, -2 %												1
<i>from</i>	-2 %, 2 %	1											

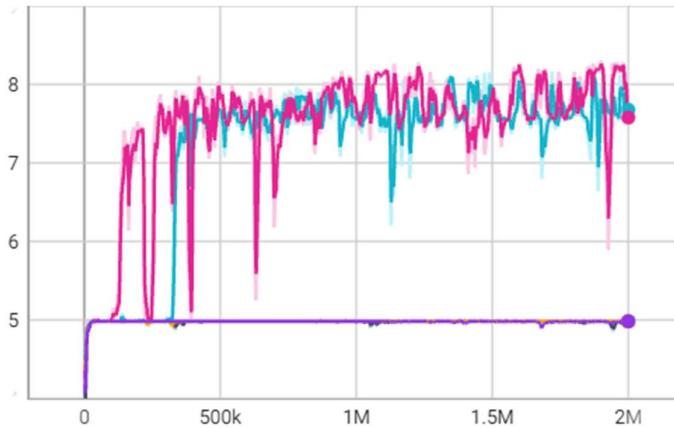
  

IE[ $r_{t+1}$ ]	Optimal
	action
1 %	Long
2 %	Long
1 %	Long
4 %	Long
-1 %	Out
5 %	Long
-2 %	Out
4 %	Long
-1 %	Out
-2 %	Out
2 %	Long
-2 %	Out

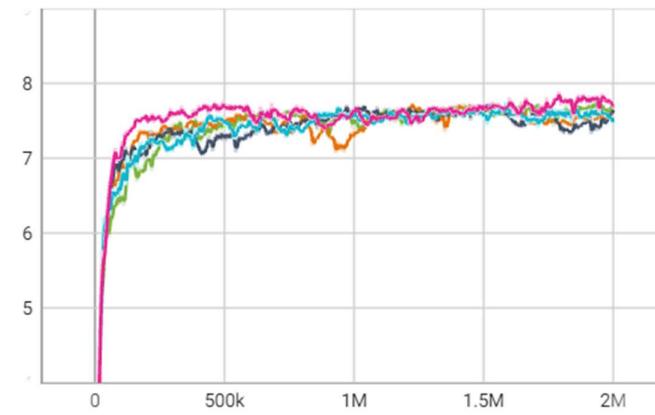
Time	1	2	3	4	5	6	7	8	9	10	11	12	Rep
Return	2 %	-2 %	1 %	2 %	1 %	4 %	-1 %	5 %	-2 %	4 %	-1 %	-2 %	2 % ...

### 3 Experiments

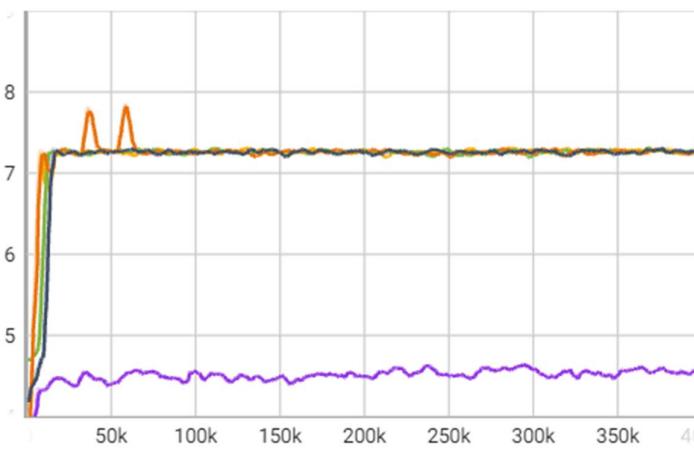
**A2C**



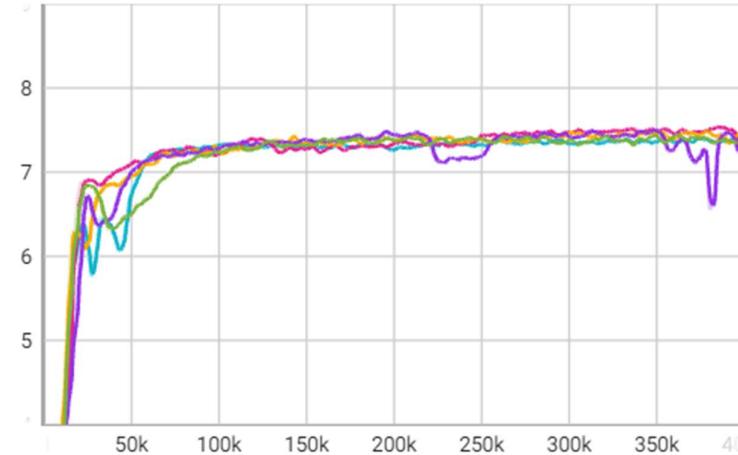
**PPO**



**DDPG**

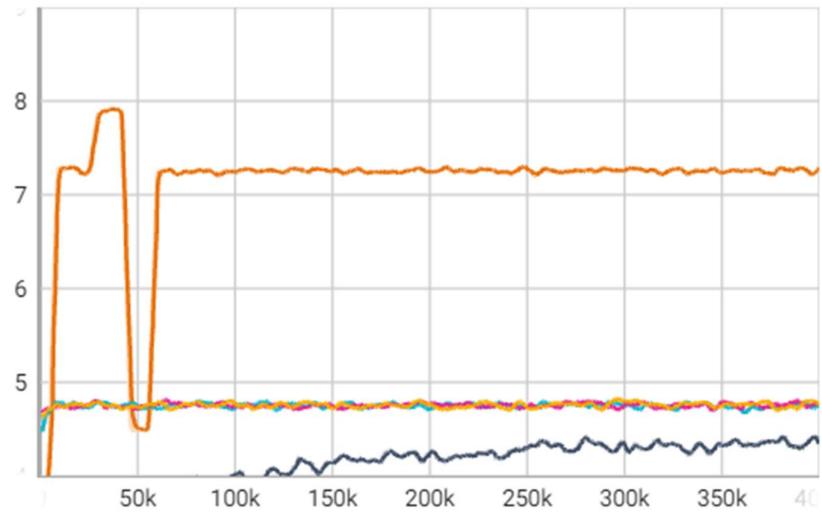


**SAC**



### 3 Experiments

#### TD3



### 3 Experiments

<b>A2C</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 54 min per run)</li><li>• 2 out of 5 runs performed well</li></ul>	<b>PPO</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 51 min per run)</li></ul> all runs performed well, but not perfect
<b>DDPG</b> <ul style="list-style-type: none"><li>• 400,000 time steps</li><li>• 5 runs (approximately 126 min per run)</li><li>• 4 of 5 runs performed well, but not perfect.</li></ul>	<b>SAC</b> <ul style="list-style-type: none"><li>• 400,000 time steps</li><li>• 5 runs (approximately 3.5 hours per run)</li><li>• All runs performed well, but not perfect.</li></ul>

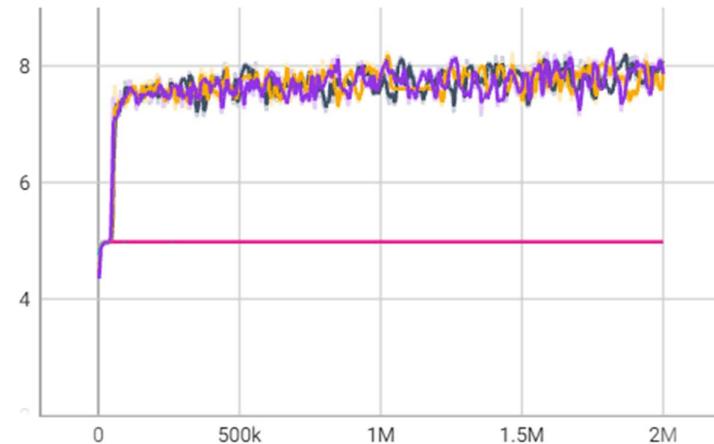
### 3 Experiments

#### TD3

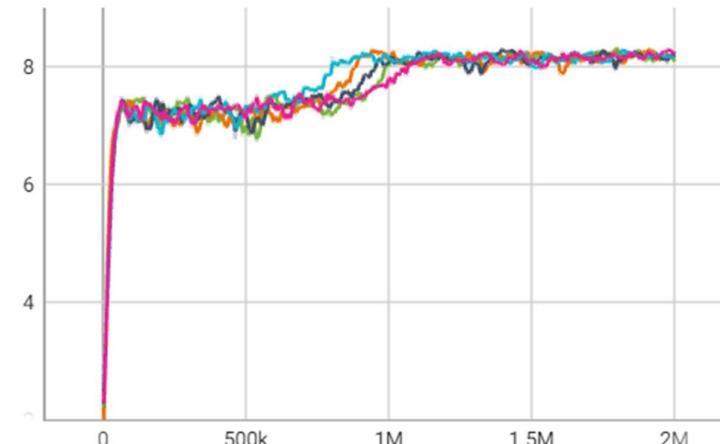
- 400,000 time steps
- 5 runs (approximately 2.7 hours per run)
- except for one run, all other runs performed badly.

### 3 Experiments

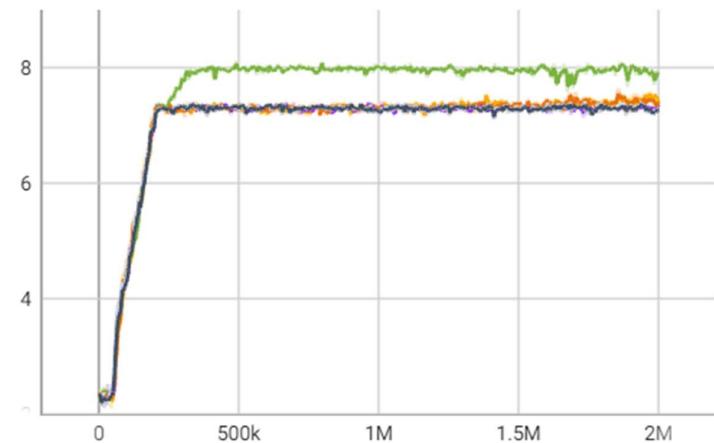
**A2C (dichotomic actions)**



**PPO (dichotomic actions)**



**DQN (dichotomic actions)**

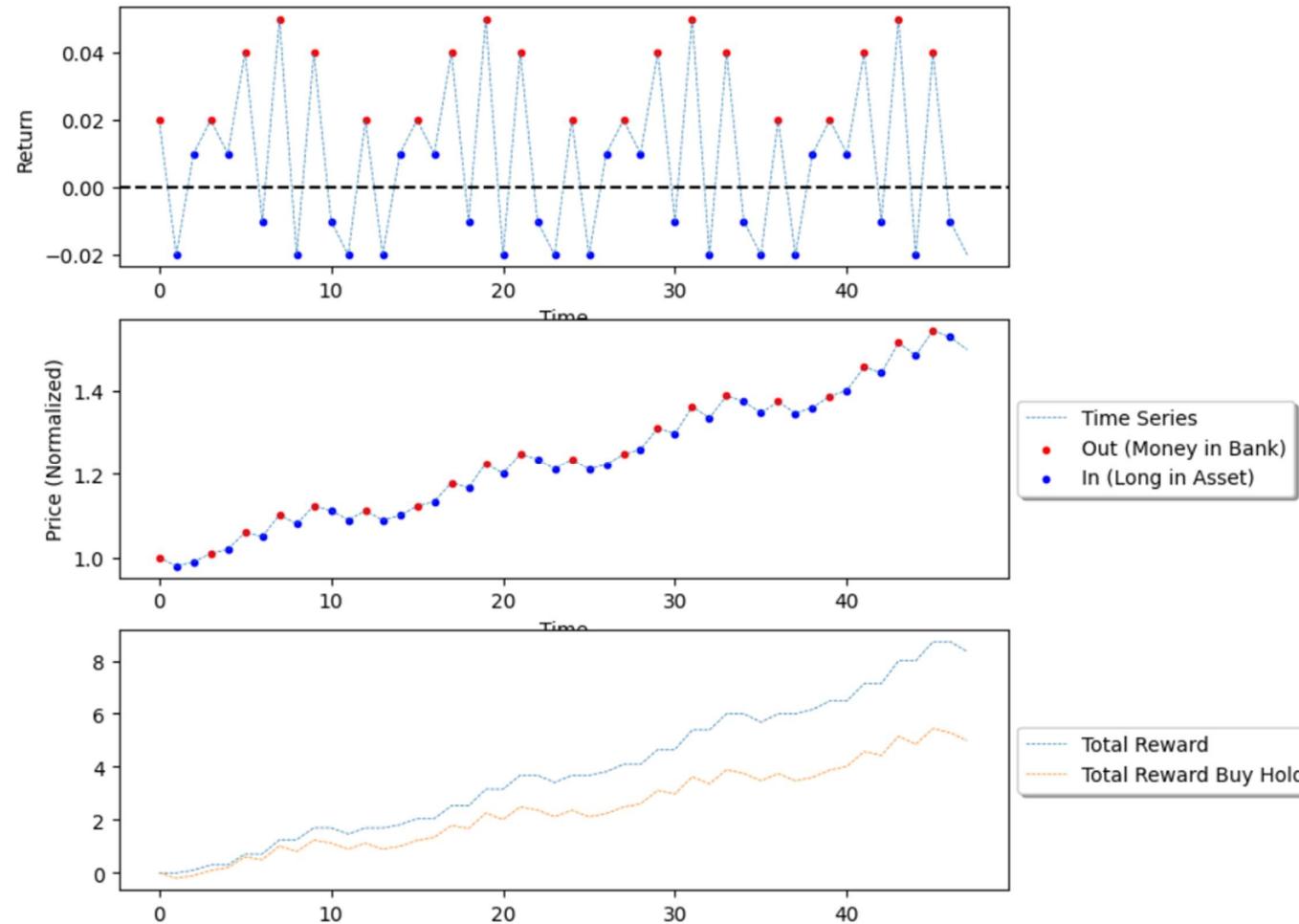


### 3 Experiments

<b>A2C (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 55 min per run)</li><li>• 3 of 5 runs performed well but not perfect (one state was wrong).</li></ul>	<b>PPO (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 52 min per run)</li><li>• all runs perform perfect.</li></ul>
<b>DQN (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 38 min per run)</li><li>• 1 of 5 worked perfect, the remaining runs performed well (one state was wrong)</li></ul>	

### 3 Experiments

#### PPO on Training Set



## 3.2 Experiment #2 – Deterministic Sine-Wave-Returns

- Time series: Deterministic sine-wave of the form:

$$r = 0.01 \cdot \sin(0.1 \cdot t) + 0.002$$

where  $t = 1, 2, \dots$

- Observation in  $t$  consists of only  $r_t$  (no information on intermediate profits, lagged returns, etc.)
- Environment #1: Continues action space in interval [0,1]
- Environment #3: Dichotomic action space: 0 or 1

### 3 Experiments

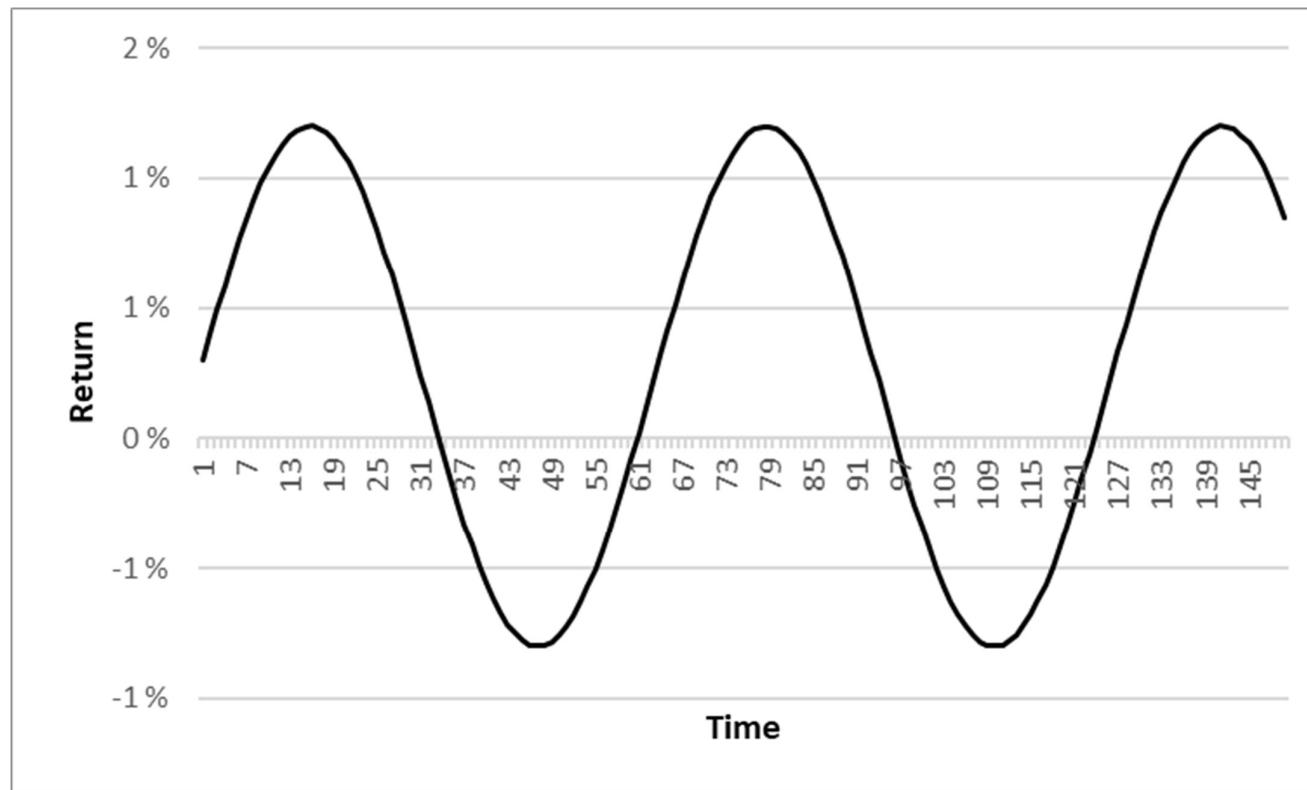


Figure 1: Sine-Wave Return (Deterministic, Discretized)

### 3 Experiments

- **Transition probabilities**, expected state returns and optimal actions (theoretically if there were 100 positive returns, followed by 0%, followed by 100 negative returns, the dataset with discrete points does not have any returns with 0%):

		<i>to</i>	<i>to</i>	<i>to</i>		
		$r > 0\%$	$r = 0\%$	$r < 0\%$	$\mathbb{E}[r_{t+1}]$	Optimal action
<i>from</i>	$r > 0\%$	99 %	1 %		>0	Long
<i>from</i>	$r = 0\%$	50 %		50 %	0	Indifferent
<i>from</i>	$r < 0\%$		1 %	99 %	<0	Out

Figure 2: Transition probabilities for (discretized) sine wave

- Training dataset consists of 135 returns (time steps).

### 3 Experiments

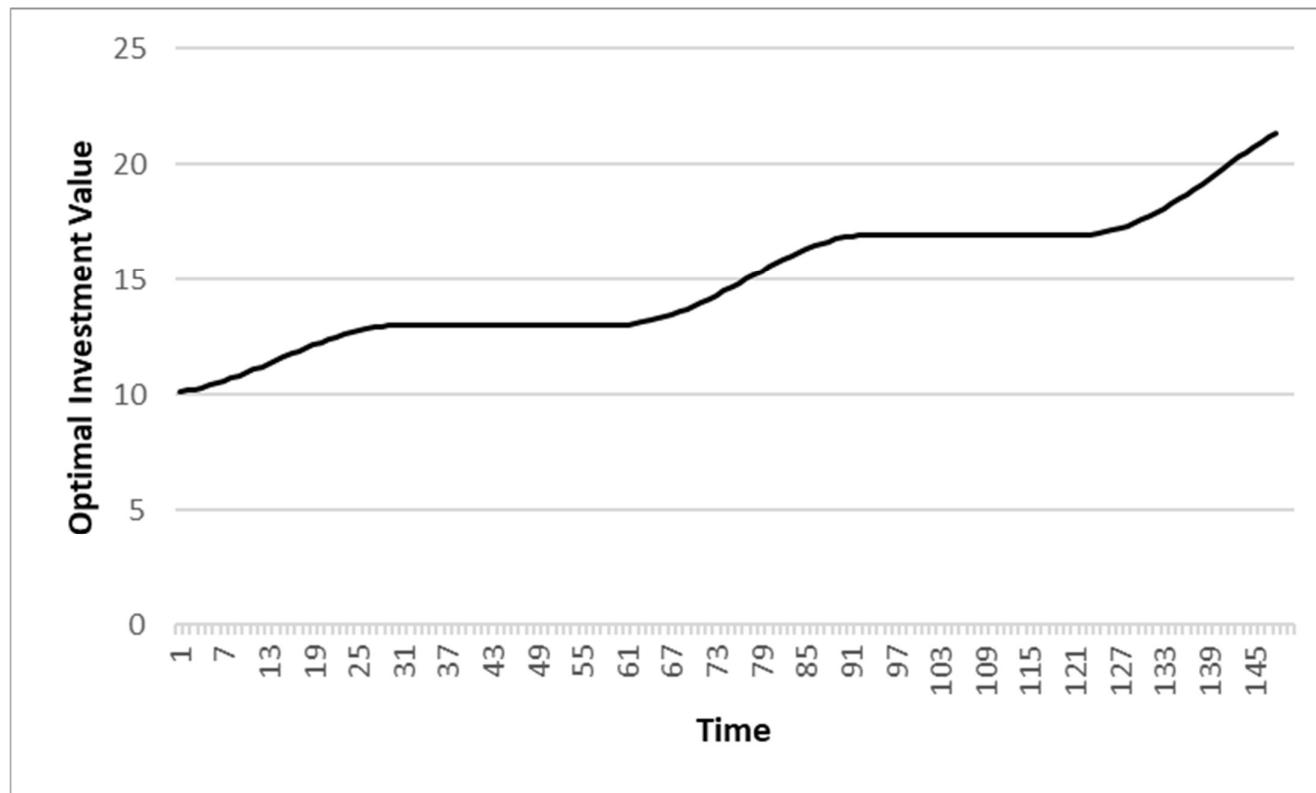
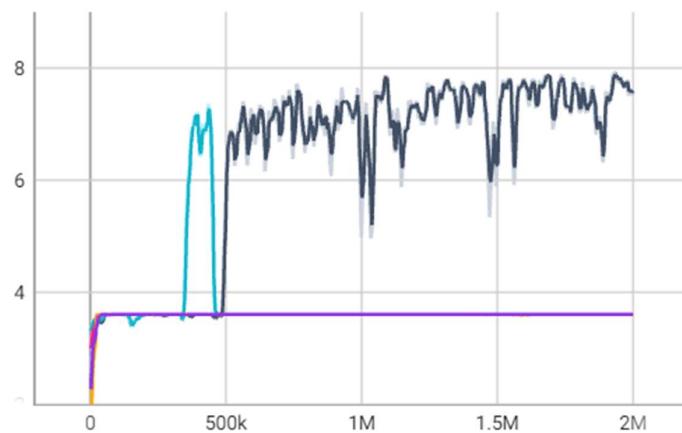


Figure 3: Optimal Trajectory of investment for sine-wave returns

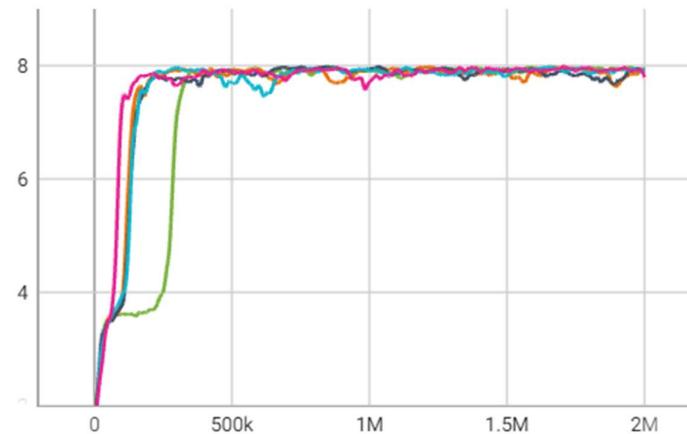
- Optimal overall return when starting with  $B_0 = 10 \Rightarrow 8.061$

### 3 Experiments

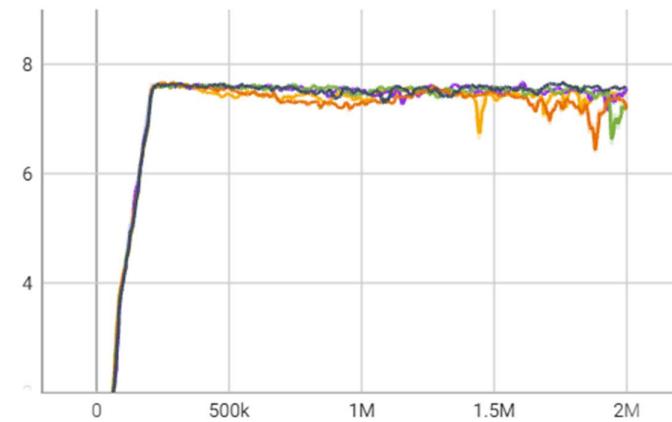
**A2C (dichotomic actions)**



**PPO (dichotomic actions)**



**DQN (dichotomic actions)**

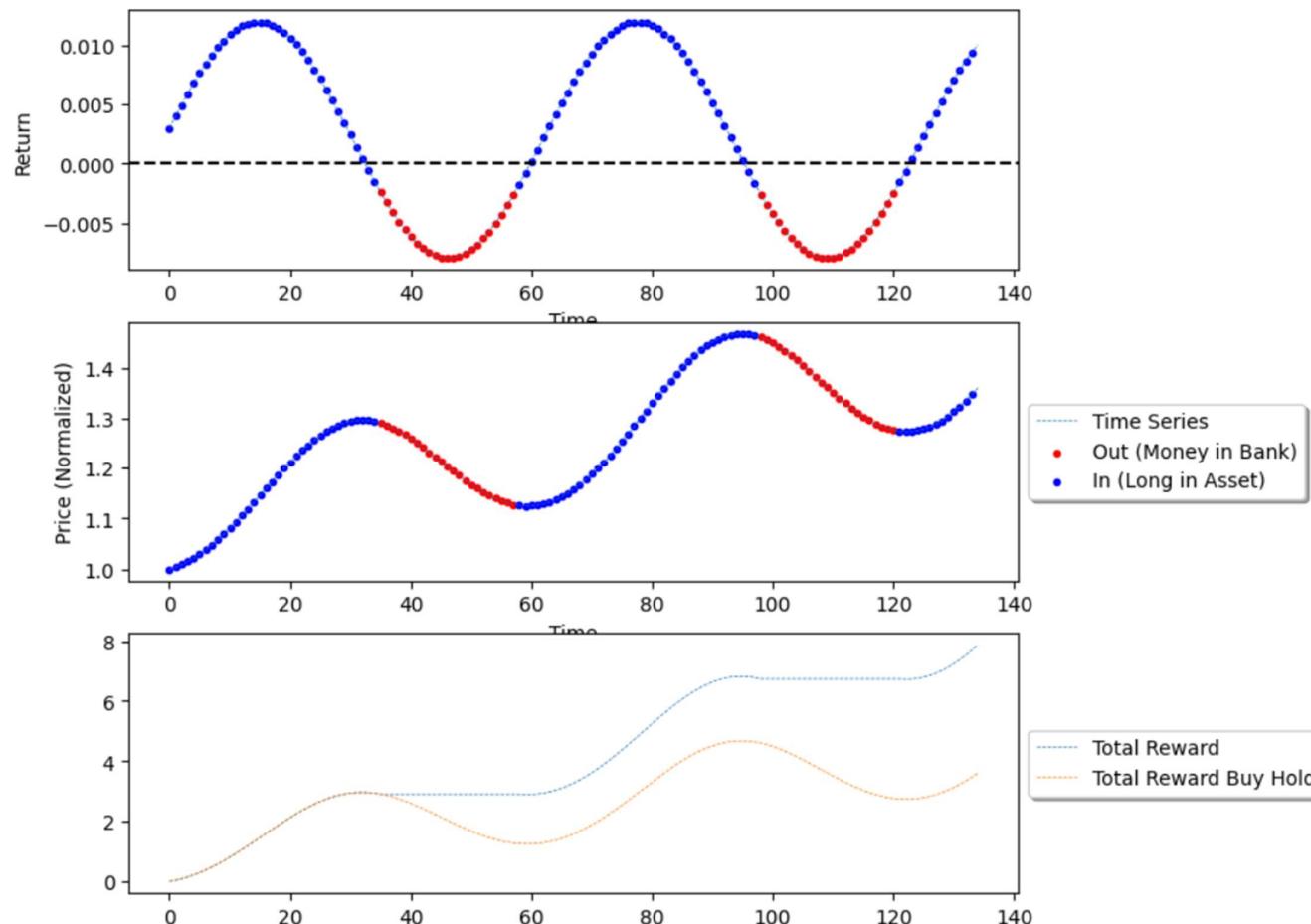


### 3 Experiments

<b>A2C (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 53 min per run)</li><li>• 1 of 5 runs performed well but not perfect (one state wrong). The best model's total return is 7.594.</li></ul>	<b>PPO (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 51 min per run)</li><li>• All runs performed well with the best model giving a total return of 7.97.</li></ul>
<b>DQN (dichotomic actions)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 33.5 min per run)</li><li>• All runs performed well, the best model achieved 7.596.</li></ul>	

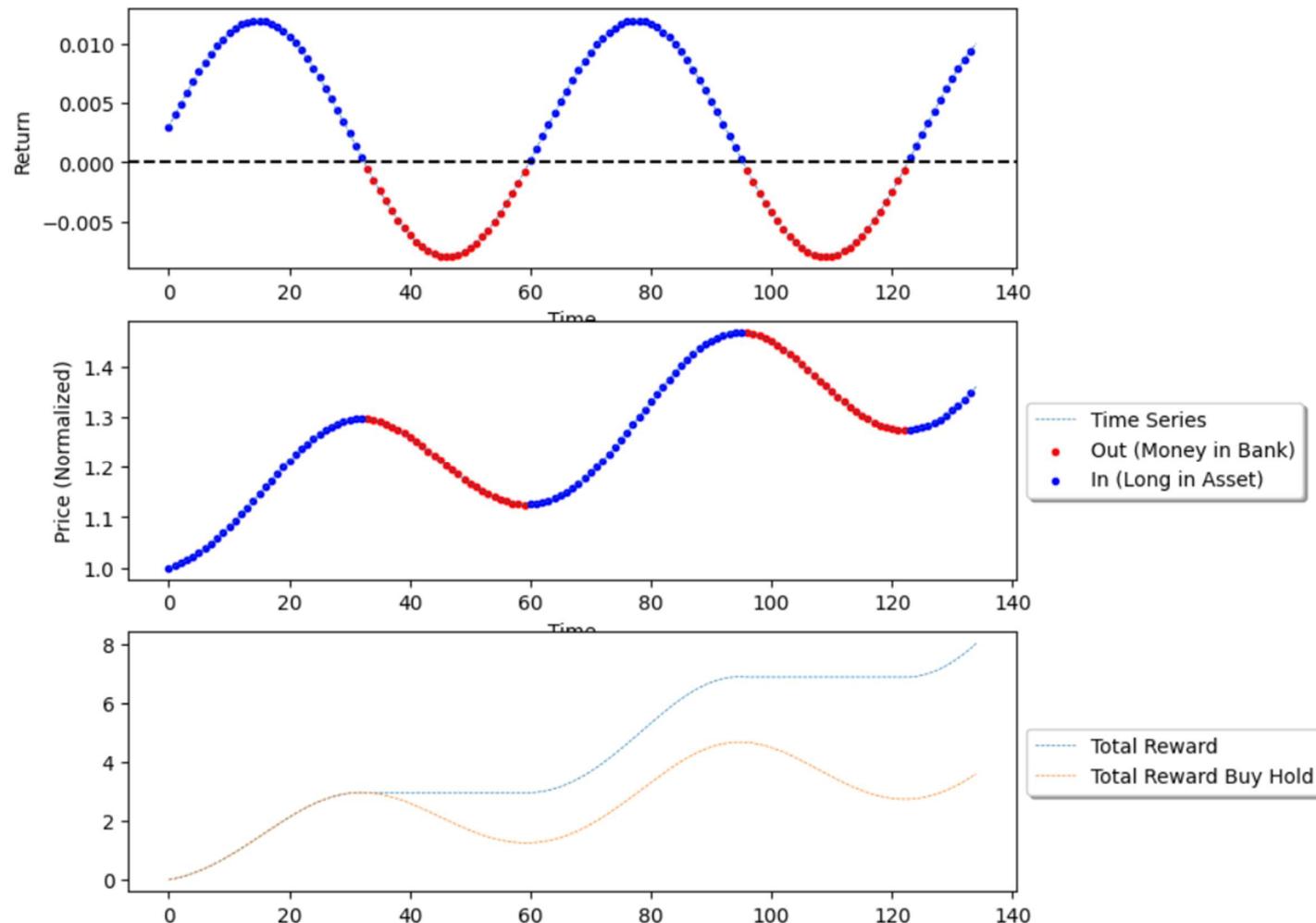
### 3 Experiments

#### DQN on Trainingsset



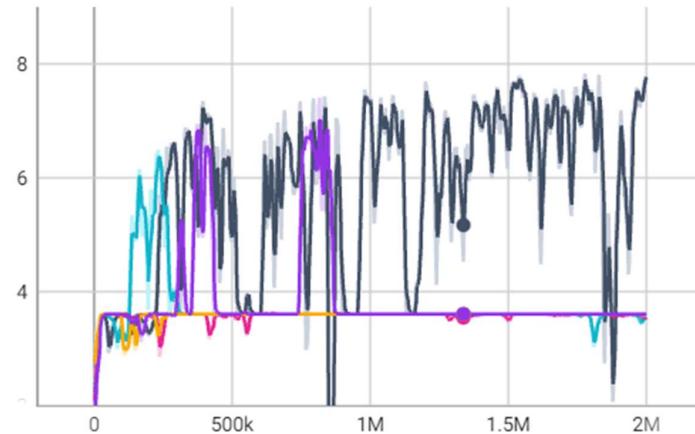
### 3 Experiments

#### PPO on Trainingsset

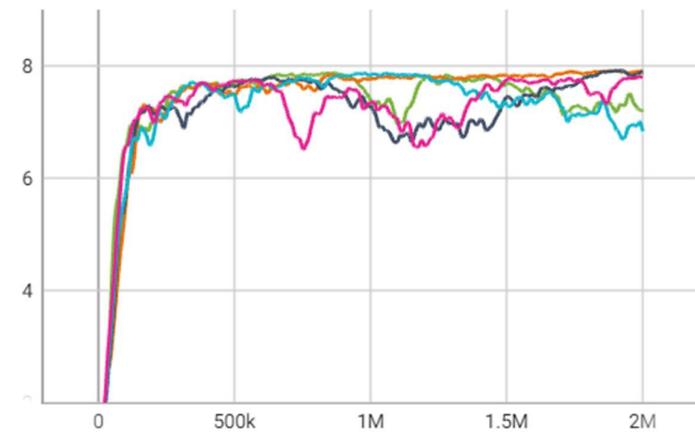


### 3 Experiments

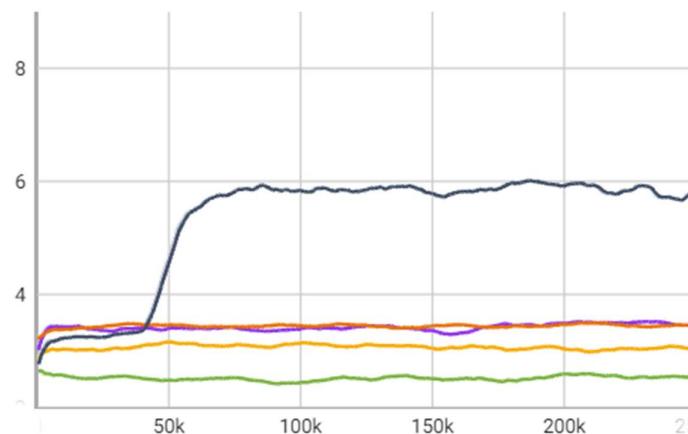
**A2C (Continuous)**



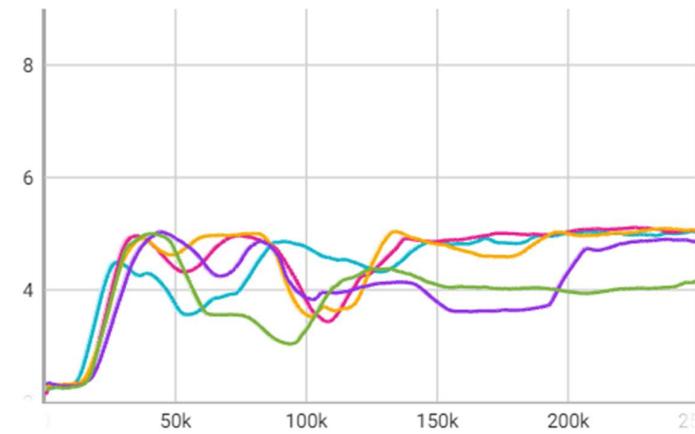
**PPO (Continuous)**



**DDPG (Continuous)**

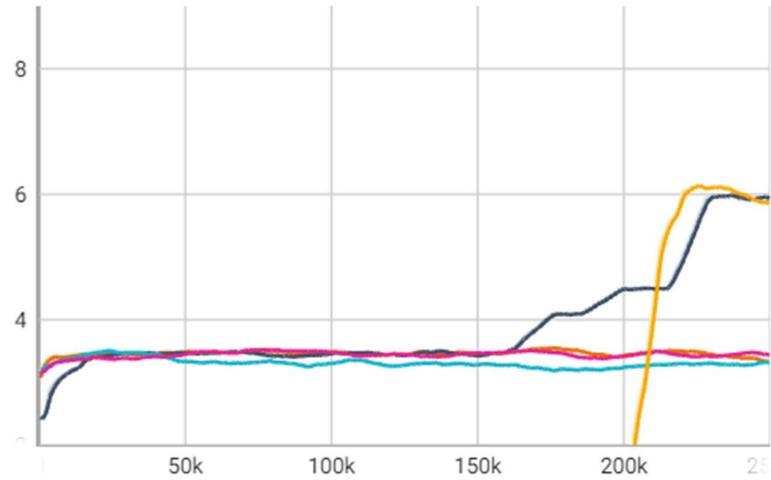


**SAC (Continuous)**



### 3 Experiments

#### TD3 (Continuous)



### 3 Experiments

<b>A2C (Continuous)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 54 min per run)</li><li>• 1 of 5 runs performed well but not perfect. The best model's total return is 7.719.</li></ul>	<b>PPO (Continuous)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 48 min per run)</li></ul> <p>5 runs performed well but not perfect. The best model's total return is 7.904.</p>
<b>DDPG (Continuous)</b> <ul style="list-style-type: none"><li>• 250,000 time steps</li><li>• 5 runs (ca. 1.18 hours per run)</li><li>• All runs not satisfying. The best model's total return is 5.819.</li></ul>	<b>SAC (Continuous)</b> <ul style="list-style-type: none"><li>• 250,000 time steps</li><li>• 5 runs (ca. 1.9 hours per run)</li><li>• All runs not very satisfying. The best model's total return is 5.069.</li></ul>

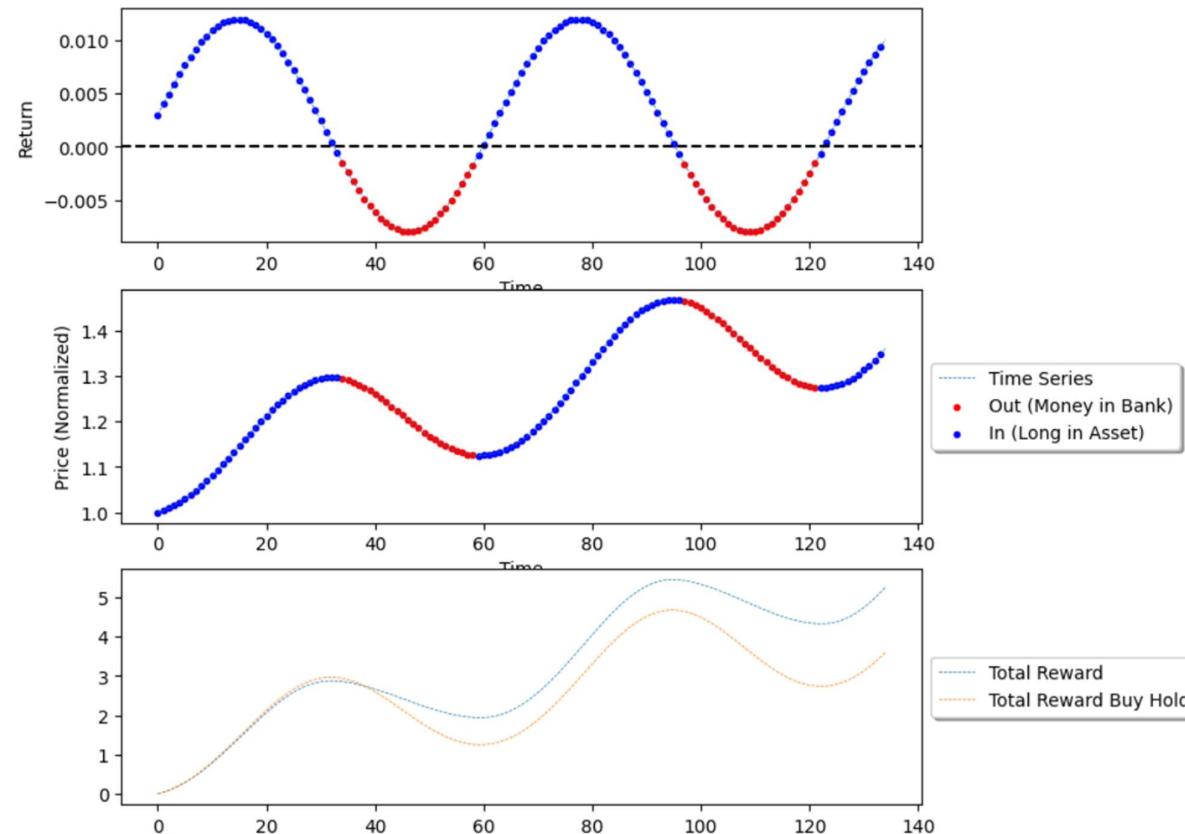
### 3 Experiments

#### **TD3 (Continuous)**

- 250,000 time steps
- 5 runs (ca. 1.3 hours per run)
- All runs not satisfying. The best model's total return is 5.912.

### 3 Experiments

#### SAC Continuous – Training Set



**Note:** Red points indicate less than 50% long; blue points indicate more than 50% long. Downward trend indicates some long position (despite indicated red) that loses value.

### 3 Experiments

#### 3.3 Experiment #3 – Geometric Brownian Motion and Regime Shifts

- Geometric Brownian Motion in Discrete Time:

$$\Delta S_t = \mu \cdot S_t \cdot \Delta t + \sigma \cdot S_t \cdot \Delta W \quad \text{where} \quad \Delta W \sim \mathcal{N}(0, \Delta t)$$

- Transition matrix for regime shifts:

		$h = 0$	$h = 1$
$h = 0$	0,99		
$h = 1$		0,97	

where 1 indicates downward trend, and 0 indicates upward trend.

$$\mu = (1 - h) \cdot \mu_{\text{Up}} + h \cdot \mu_{\text{Down}}$$

$$\sigma = (1 - h) \cdot \sigma_{\text{Up}} + h \cdot \sigma_{\text{Down}}$$

- Theoretically shifts are not predictable from the  $r_t$ , because shifts do not depend on it.

### 3 Experiments

- Shifts may be identified in the hindsight, when the average return becomes visible.

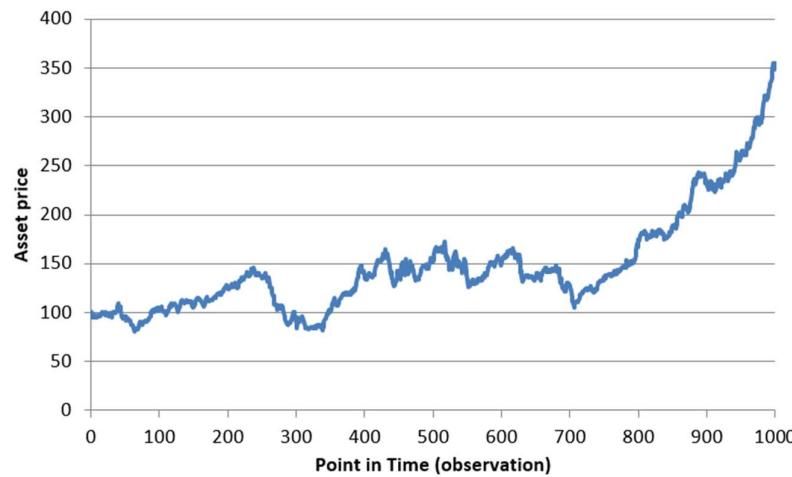


Figure 4: Observed trajectory for GBM with regime shifts

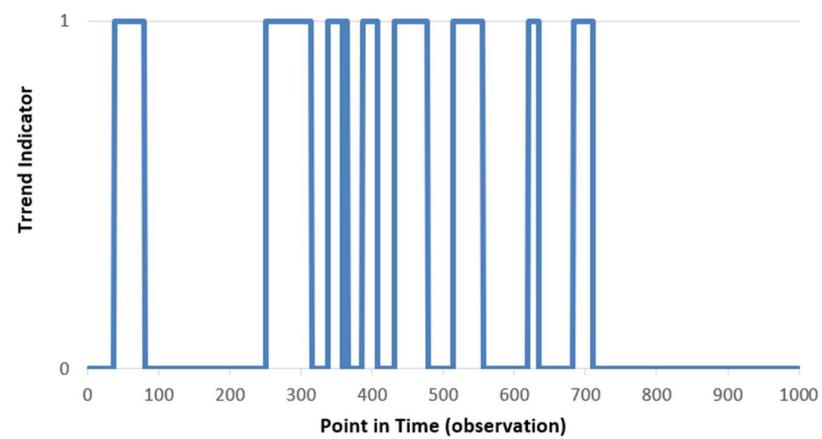
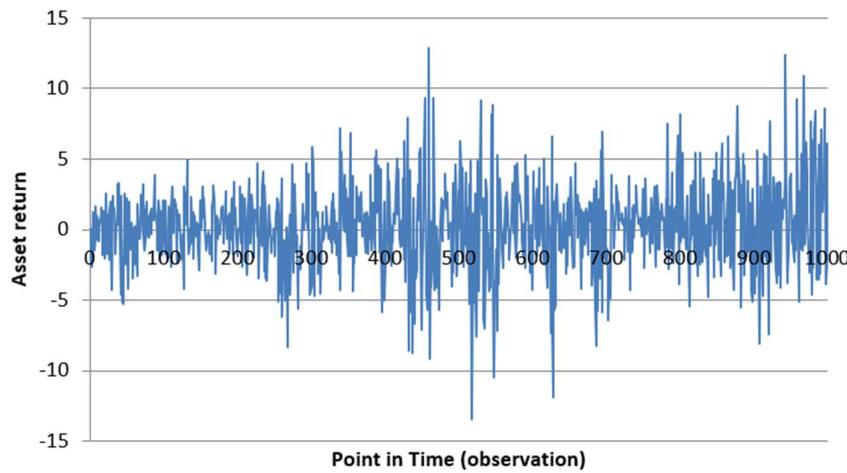


Figure 5: Observed Regime Shifts

### 3 Experiments



**Figure 6: Observed returns for GBM with regime shifts**

### 3 Experiments

**State space =  $\{r_t\}$**

- The trend cannot be spotted (identified) with  $r_t$  as the only information.
- $r_{t+1}$  cannot be predicted.
- The optimal decision of reward maximizer: always in if whole series trends upwards, always out if whole series trends downwards.

**State space =  $\{r_t, h_t\}$**

- Perfect information concerning the trend.
- $r_{t+1}$  cannot be predicted.
- The optimal decision of reward maximizer: always in if  $h = 0$ , always out if  $h = 1$ .

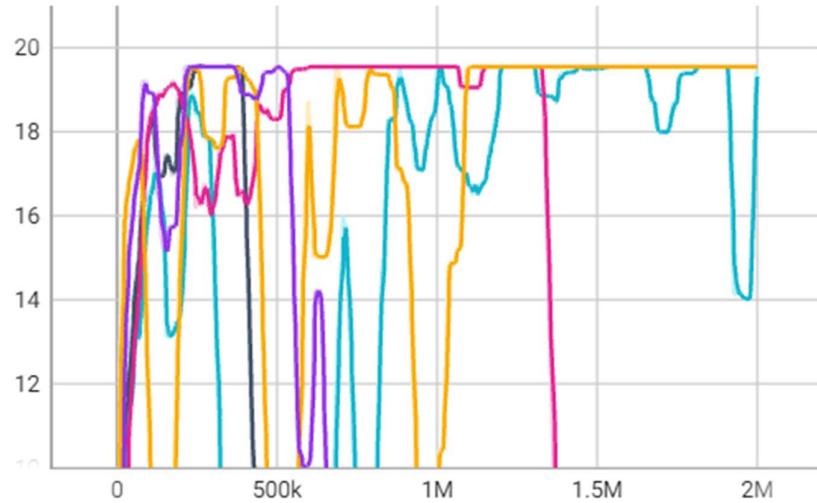
### 3 Experiments

**State space =  $\{r_t, r_{t-1}, r_{t-2}\}$**

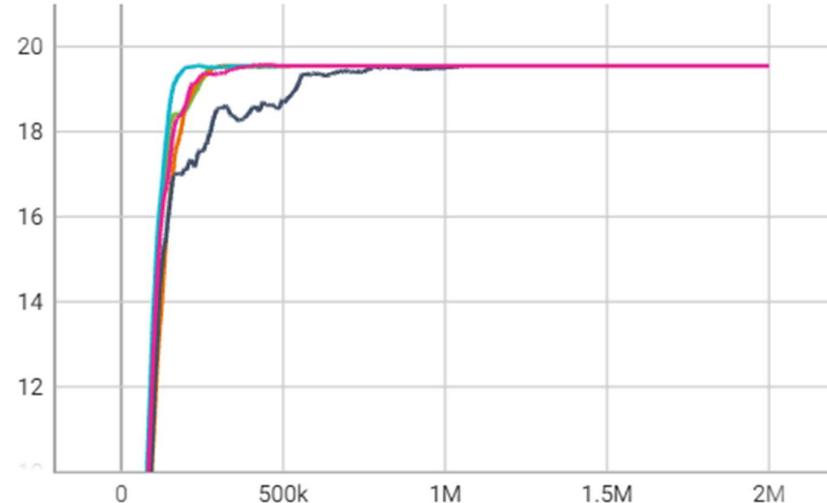
- The trend cannot be predicted (in  $t - 1$ ) nor will it be identified immediately (in  $t$ ).
- However, the sequence of lagged returns may give some information on the trend.
- This information will depend on length of lags and  $\sigma$ .
- Expected behavior: more in than out if  $h = 0$ , more out than in if  $h = 1$ .

### 3 Experiments

**A2C (Continuous)**



**PPO (Continuous)**

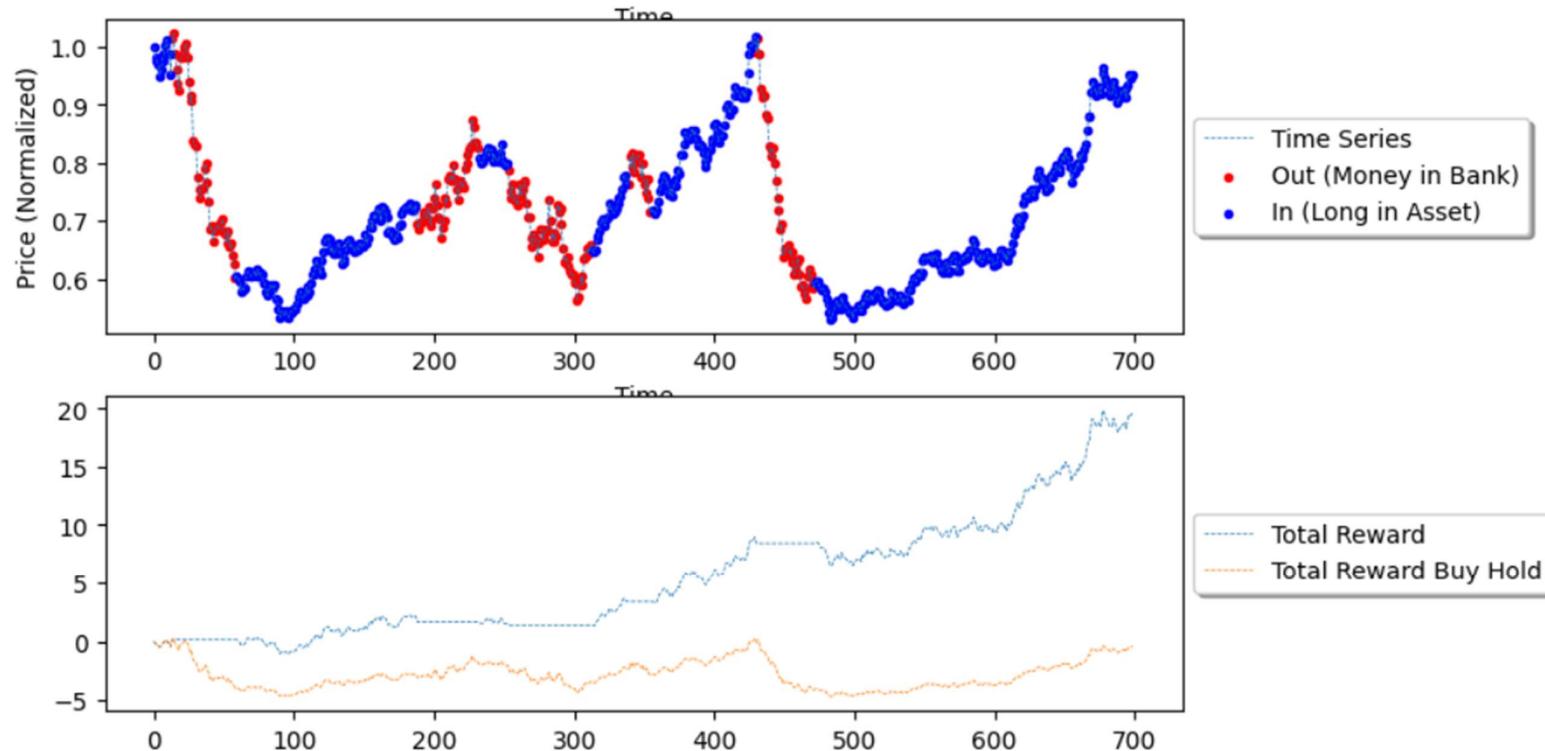


### 3 Experiments

A2C (Continuous)	PPO (Continuous)
<ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 52 minutes per run)</li><li>• Very unstable, but all five models performed well at some point.</li><li>• All runs not satisfying. The best model's total return is 19.55.</li></ul>	<ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (approximately 52 minutes per run)</li><li>• All runs converged to a total return of 19.55.</li></ul>

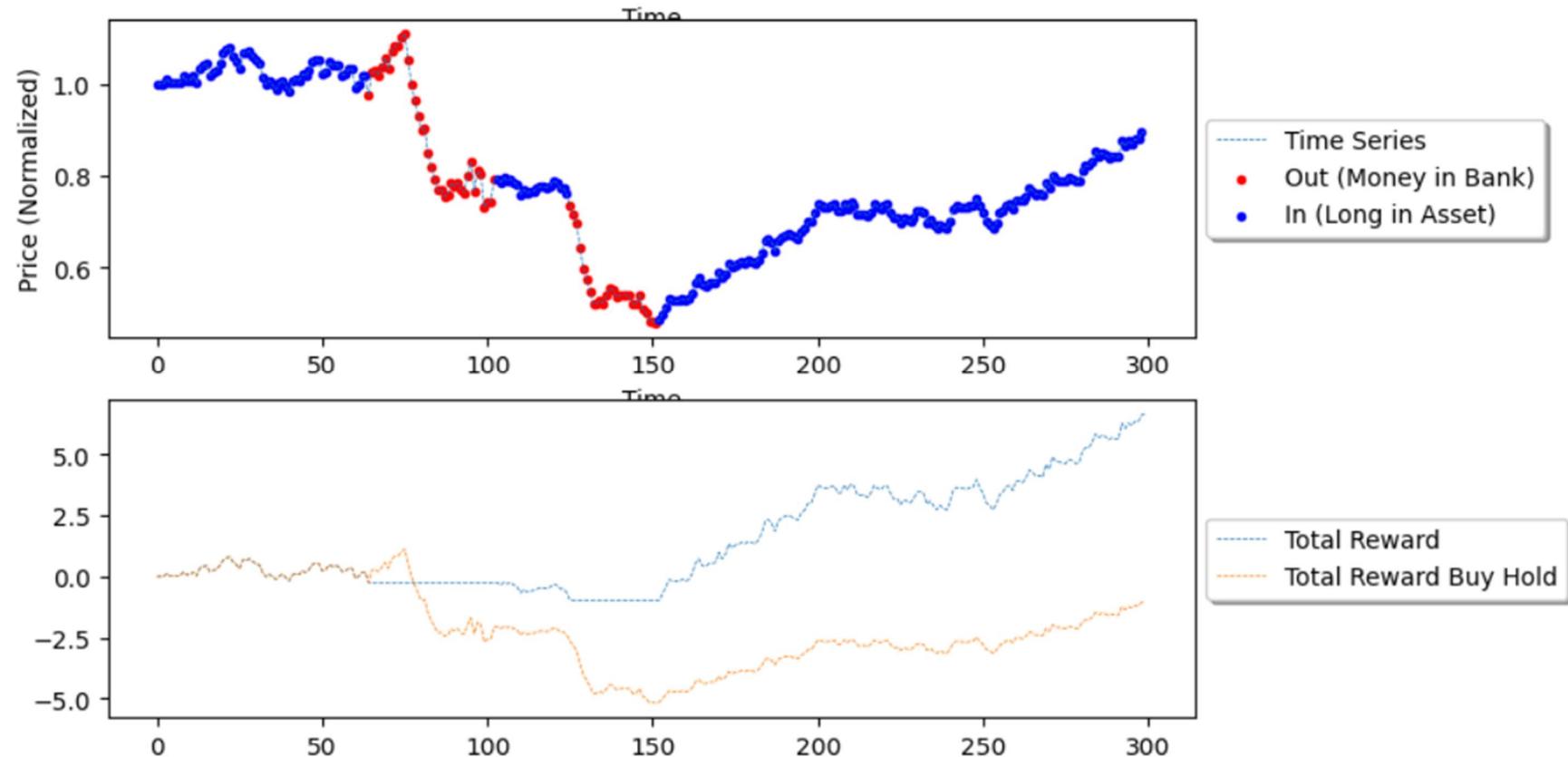
### 3 Experiments

#### Best PPO on Training Set



### 3 Experiments

#### Best PPO – Test Set



.

### 3 Experiments

#### **DDPG (Continuous)**

- 250,000 time steps
- 5 runs (approximately 1,32 hours per run)
- All runs performed badly, no total return over 11.11

#### **DDPG (Continuous)**

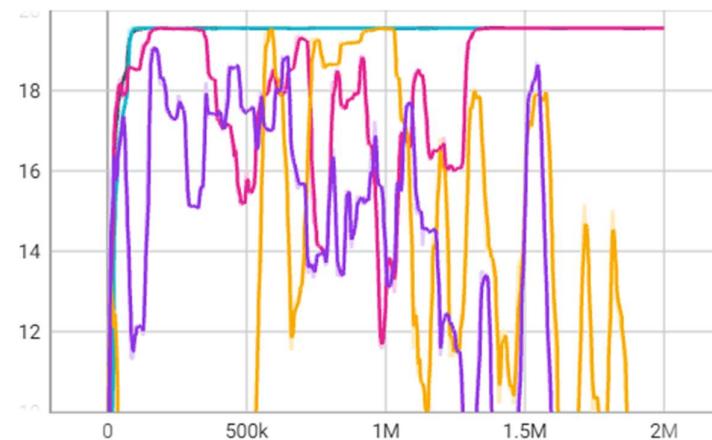
- 250,000 time steps
- 5 runs (approximately 2,1 hours per run)
- All runs performed badly, all total returns below 9.

#### **TD3 (Continuous)**

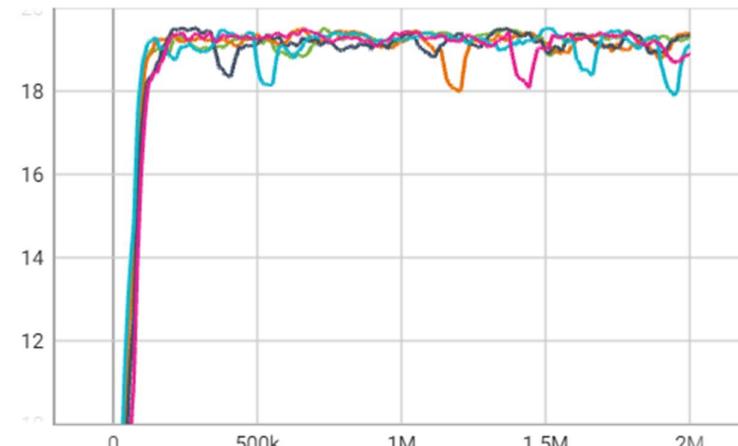
- 250,000 time steps
- 5 runs (approximately 1.47 hours per run)
- All runs performed badly, all total returns below 11.

### 3 Experiments

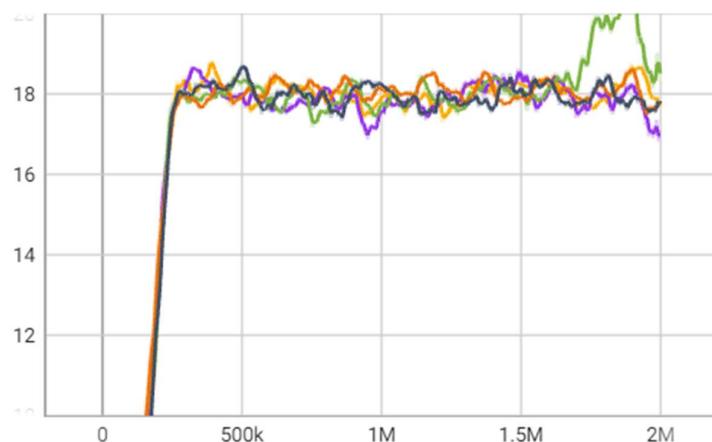
**A2C (binary)**



**PPO (binary)**



**DQN (binary)**

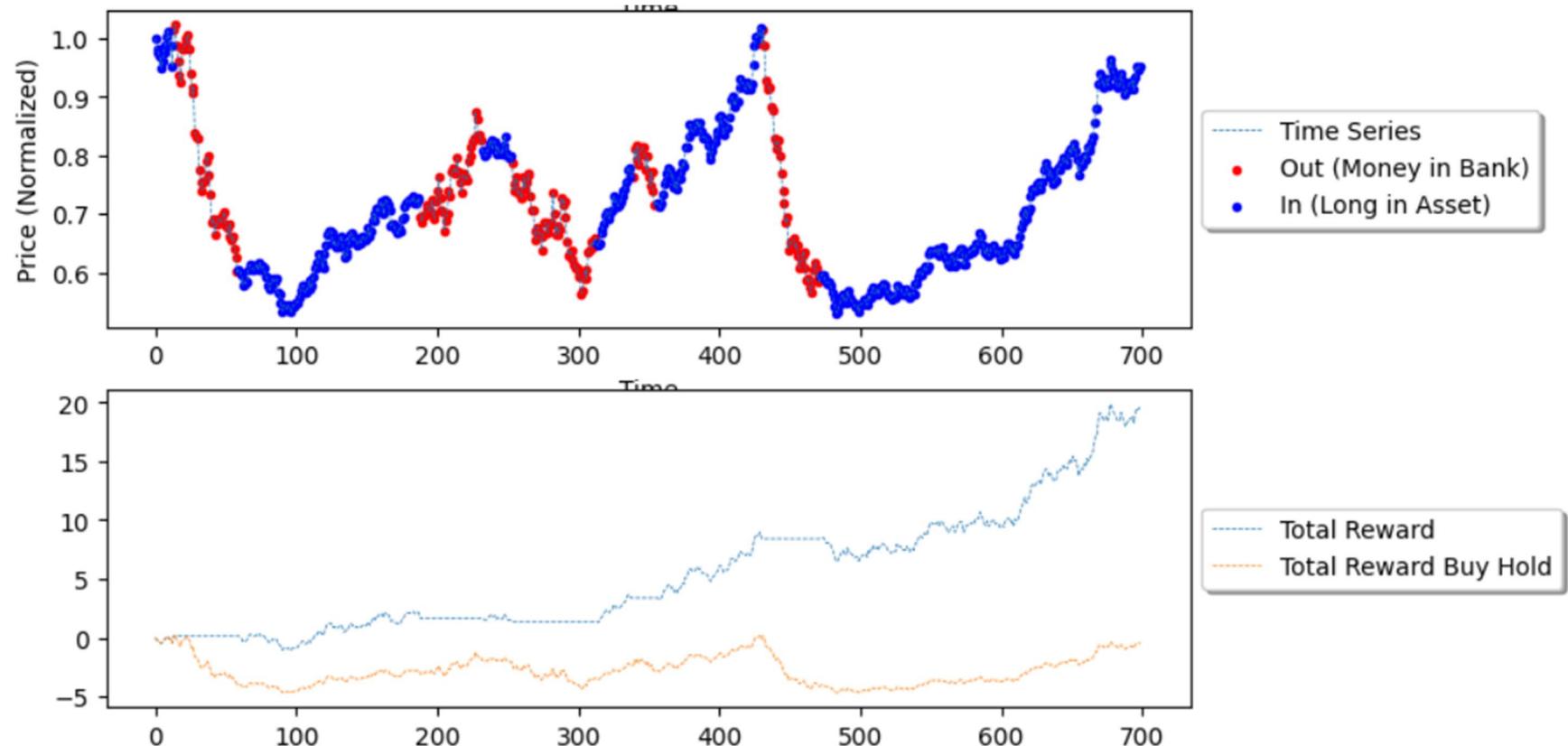


### 3 Experiments

<b>A2C (binary)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 53 minutes per run)</li><li>• Unstable, four of five models reached 19:55 at some point)</li></ul>	<b>PPO (binary)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 51 minutes per run)</li><li>• All runs performed quite well, best model in the end 19.37 (at some point 19.54)</li></ul>
<b>DQN (Discrete)</b> <ul style="list-style-type: none"><li>• 2,000,000 time steps</li><li>• 5 runs (ca. 34 minutes per run)</li><li>• All runs performed quite well, best model at some point 20.53; at end point 18.3)</li></ul>	

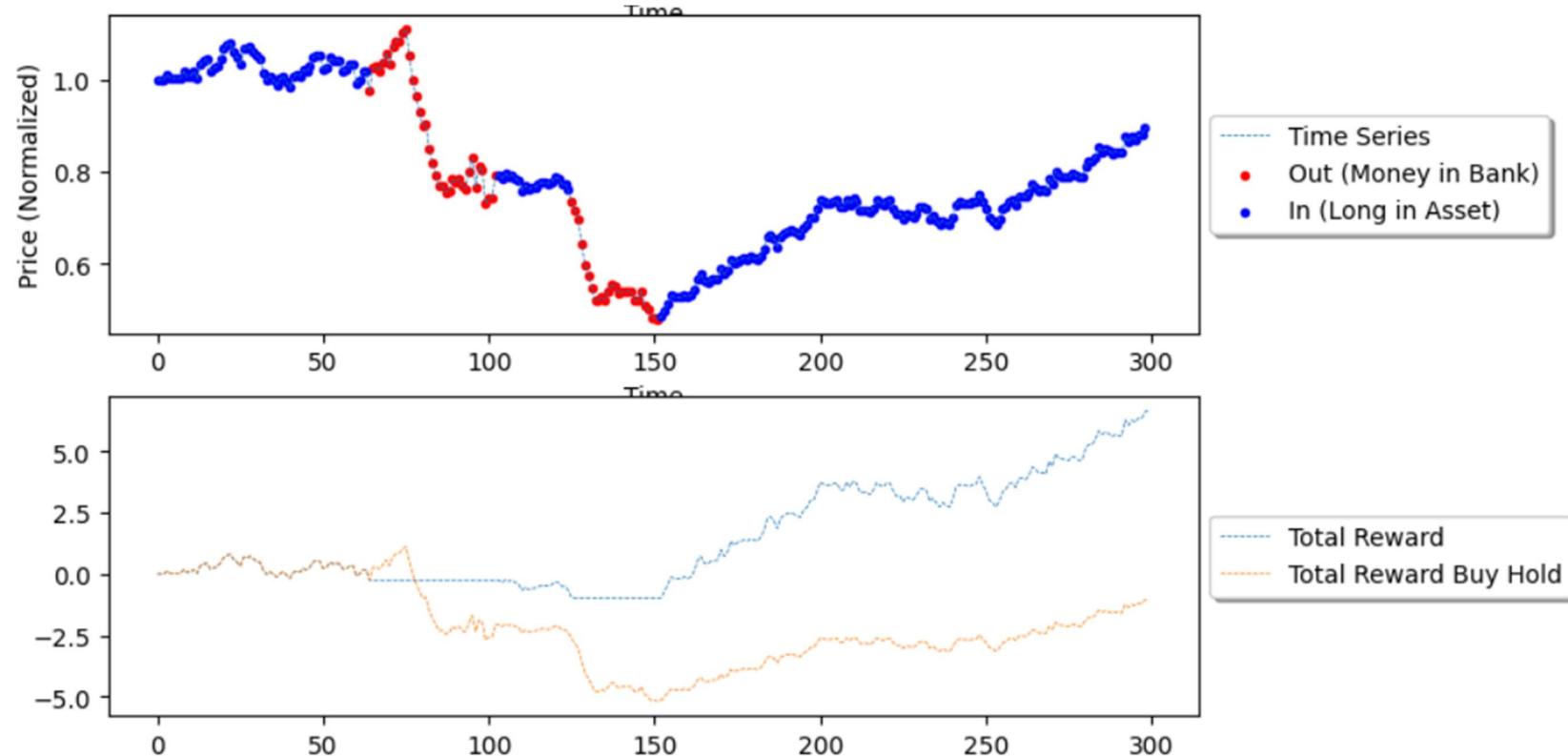
### 3 Experiments

#### Best PPO on Training Set



### 3 Experiments

#### Best PPO – Test Set



### 3 Experiments

<b>Obs (all)</b>	<b>Obs (test)</b>	<b>State</b>	<b>Training</b>	<b>Test</b>
0-13		0	Nice	
14-59		1	Nice	
60-188		0	Nice	
189-232		1	Nice	
233-253		0	Nice	
254-313		1	Nice	
314-338		0	Nice	
339-356		1	Nice	
357-430		0	Nice	
431-472		1	Nice	
473-764	xx-64	0	Nice	Nice
765-803	65-103	1		Nice
804-825	104-125	0		Nice
826-852	126-152	1		Nice
853-1000	153-300	0		Nice

### 4 Take-aways, To-Do's and Future research

#### Take-aways:

- Some RL methods in SB3 deliver “bad” performance even for simple deterministic time series.
- PPO seems to perform best for a variety of time series, and state spaces.
- The action space (discrete versus continuous) affects the solution even if the underlying problem and its solution should be the same for both action spaces.
- Problem is quickly overfitted: Problematic with respect to the features that can be added to the state space.

### To Do:

- Train/Test with more data.
- Test on more test sets.
- Implement on-training cross validation (not really available by default in SB3)
- Implement early stopping, best model selection (A2C).
- Implement the risk-return trade-off