

RL FINAL PROJECT PRESENTATION

YASH ADIVAREKAR - 2021101008

YASH KAWADE - 2021101032

MADHAV TANK - 2021101108

MANAV SHAH - 2021101090

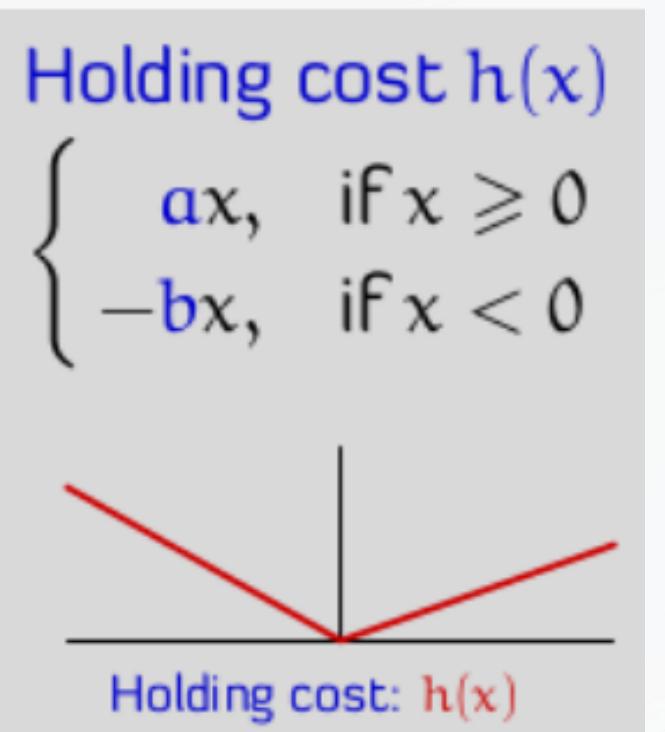
NIKUNJ GARG - 2021101021



INVENTORY MANAGEMENT PROBLEM

- Retail stores stockpile products in warehouses to meet the random demand. Additional stocks are procured at regular intervals.
- Let X_t denote the amount of stock before the t -th procurement.
- At time t , the store may procure an additional stock U_t ($\leq U$) units for price p per unit. Thus the total procurement cost is pU_t .

- The random demand W_t is i.i.d. with distribution P_w . P_w is uniform[0,10].
- The stock available at the next time is $X(t+1) = X_t + U_t - W_t$, where a negative stock denotes backlogged demand.
- Per-stage cost is $c(X_t, U_t) = h(X_t) + pU_t$. The holding cost for the stock is given by $h(x)$ where a is the per-unit storage cost and b is the per-unit backlog cost.
- We want to find the optimal inventory control strategy to minimize the expected total cost over a finite horizon.



State Space

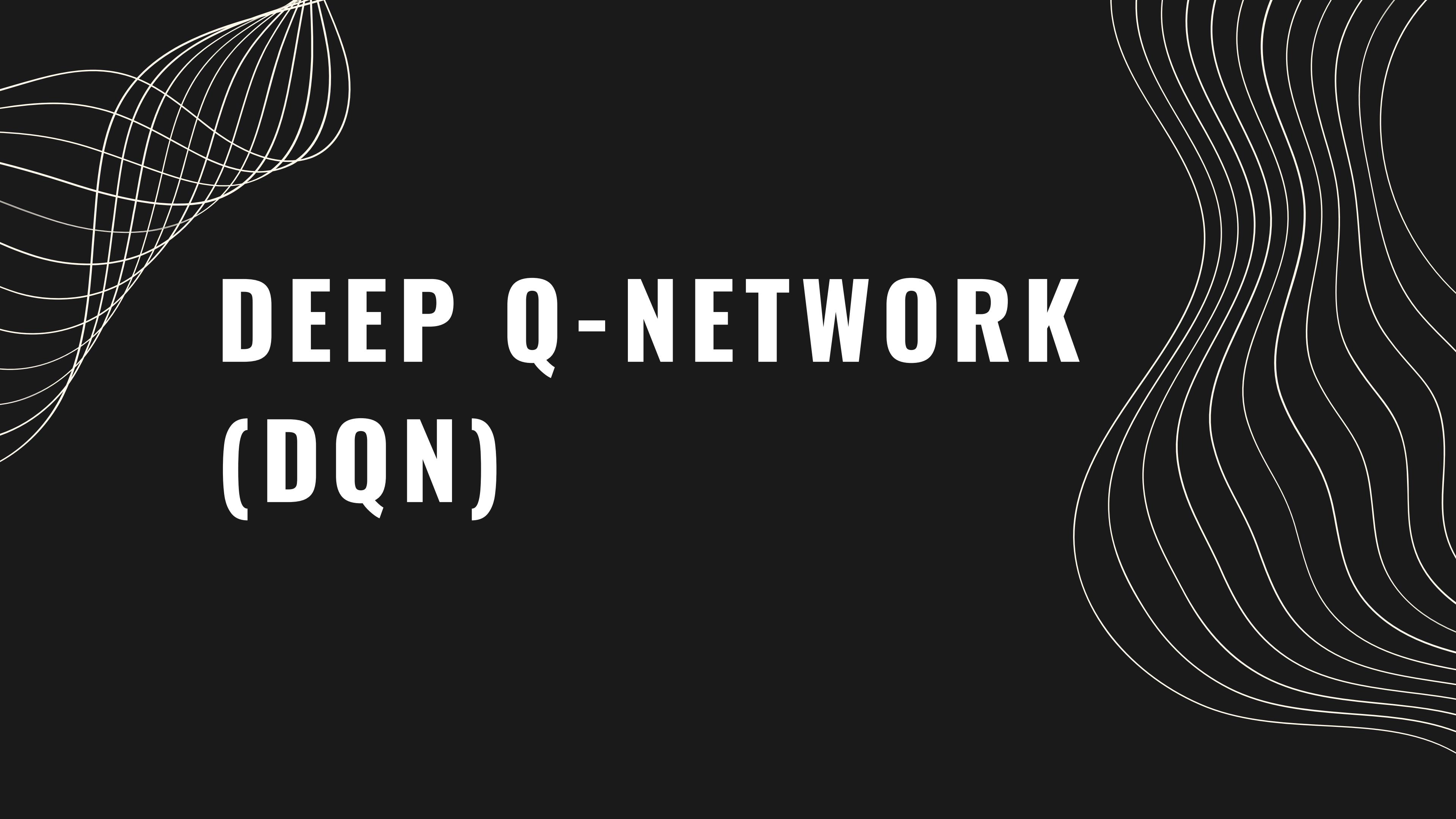
Our states are X 's (the amount of inventory). Our state's range is X 's range.

$$-10 \leq X \leq 10$$

Actions

The amount of inventory produced in that day is taken as the action. $U(t)$ is our action on t^{th} day.

$$0 \leq U \leq 10$$



DEEP Q-NETWORK (DQN)

Pseudo code

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}

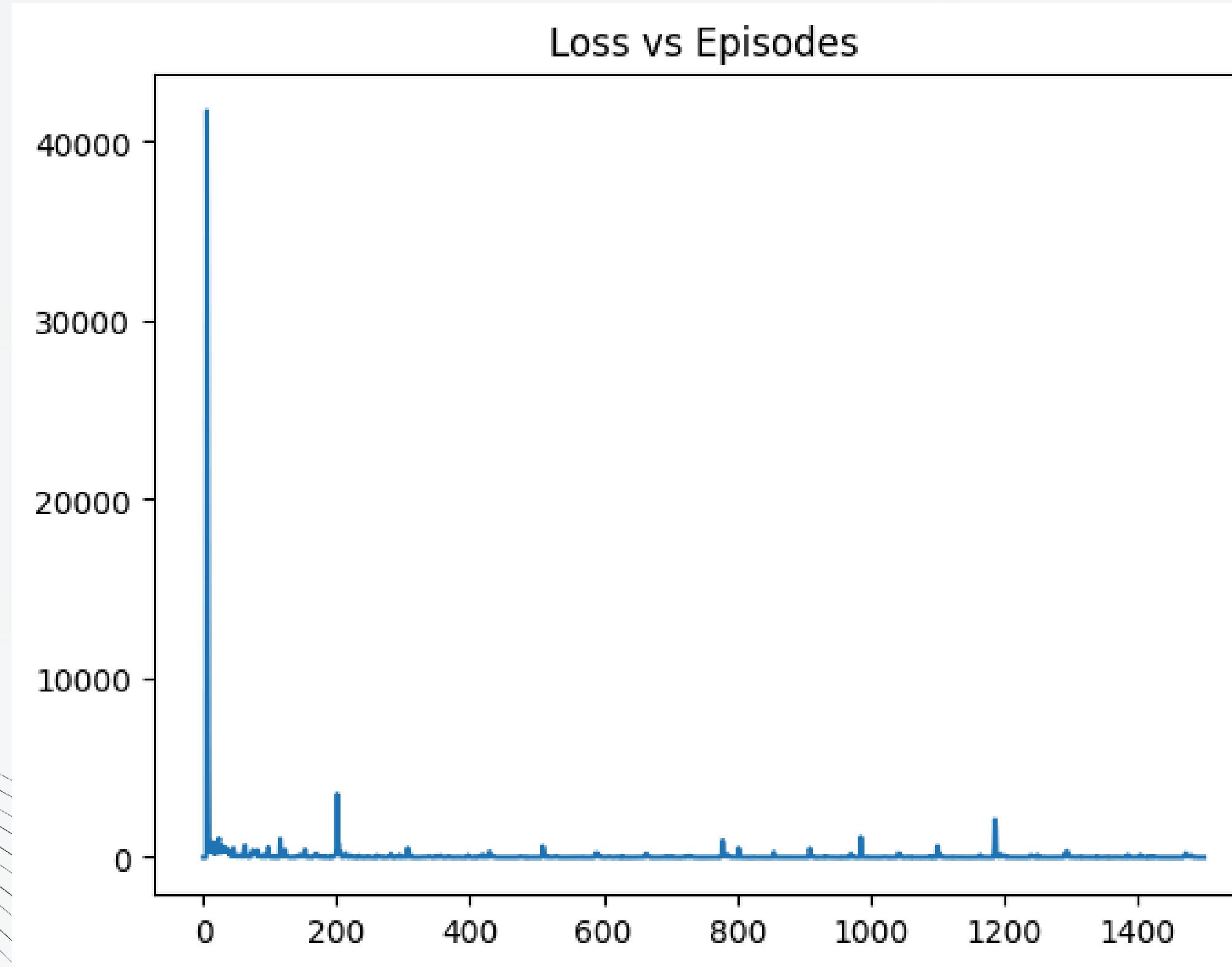
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$

end for

end for

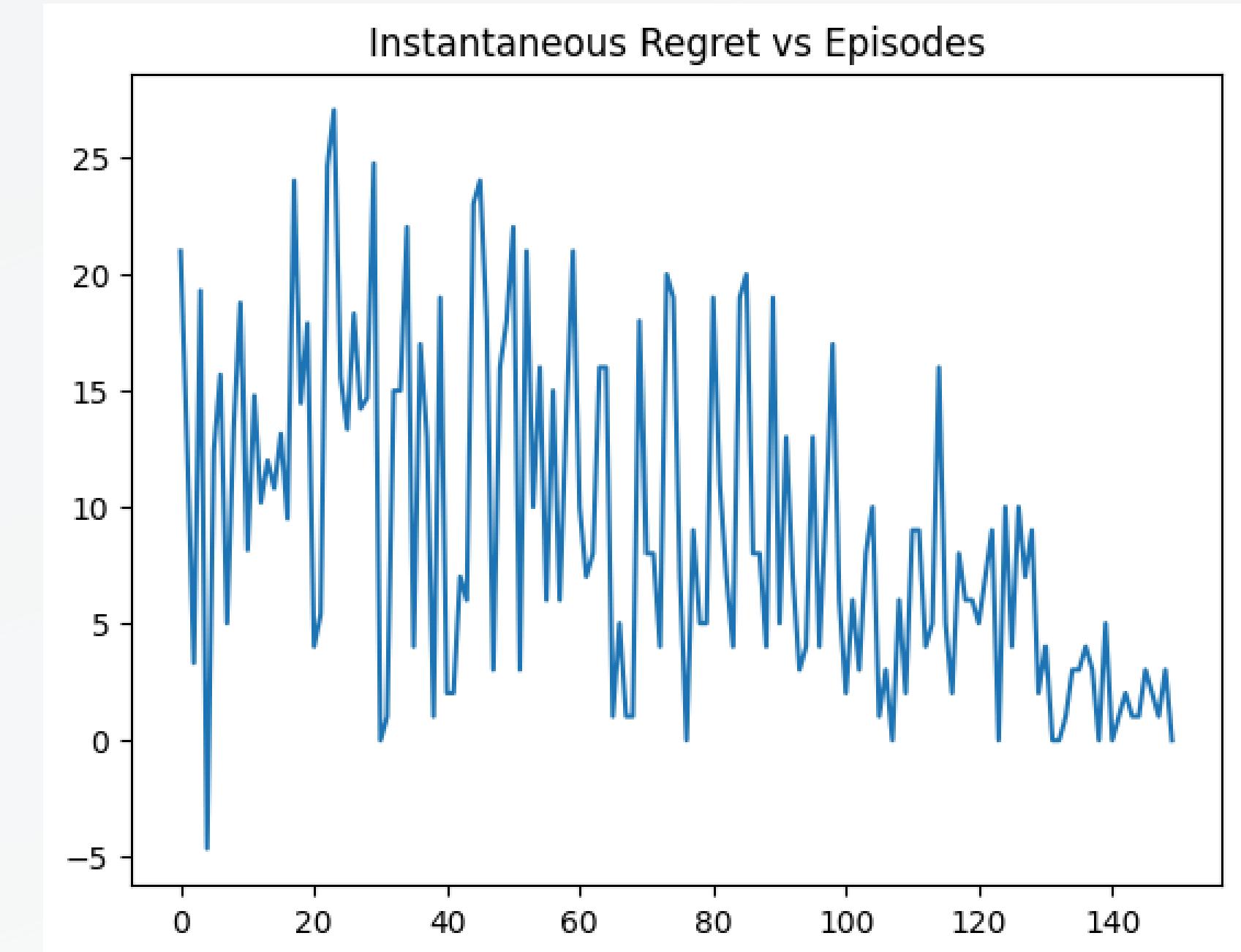
DQN Loss plot

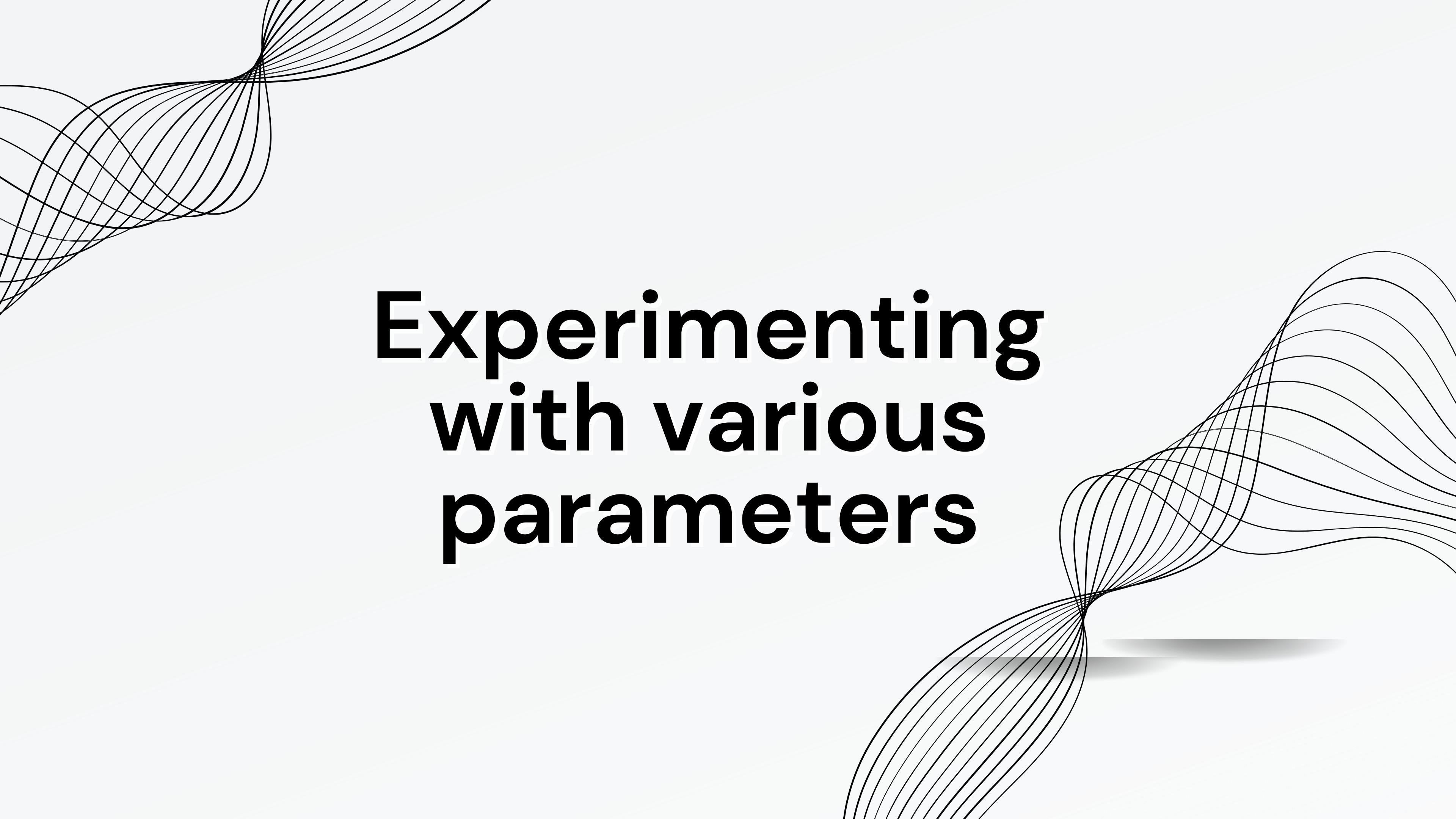


Cumulative Regret



Instantaneous Regret





Experimenting with various parameters

1

The DQN algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 3
 - Backlog Cost (b): 3
 - Purchase Cost per Unit (p): 0
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The DQN algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 50
 - Backlog Cost (b): 5
 - Purchase Cost per Unit (p): 15
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

States from -5 to 5

	-5	-4	-3	-2	-1	0	1	2	3	4	5
0	0	1	2	1	0	0	0	0	0	0	0
1	1	0	5	5	4	3	2	1	0	0	0
2	2	5	5	4	3	2	1	0	0	0	0
3	3	5	4	3	2	1	0	0	0	0	0
4	4	4	3	2	1	0	0	0	0	0	0
w	w	3	2	1	0	0	0	0	0	0	0

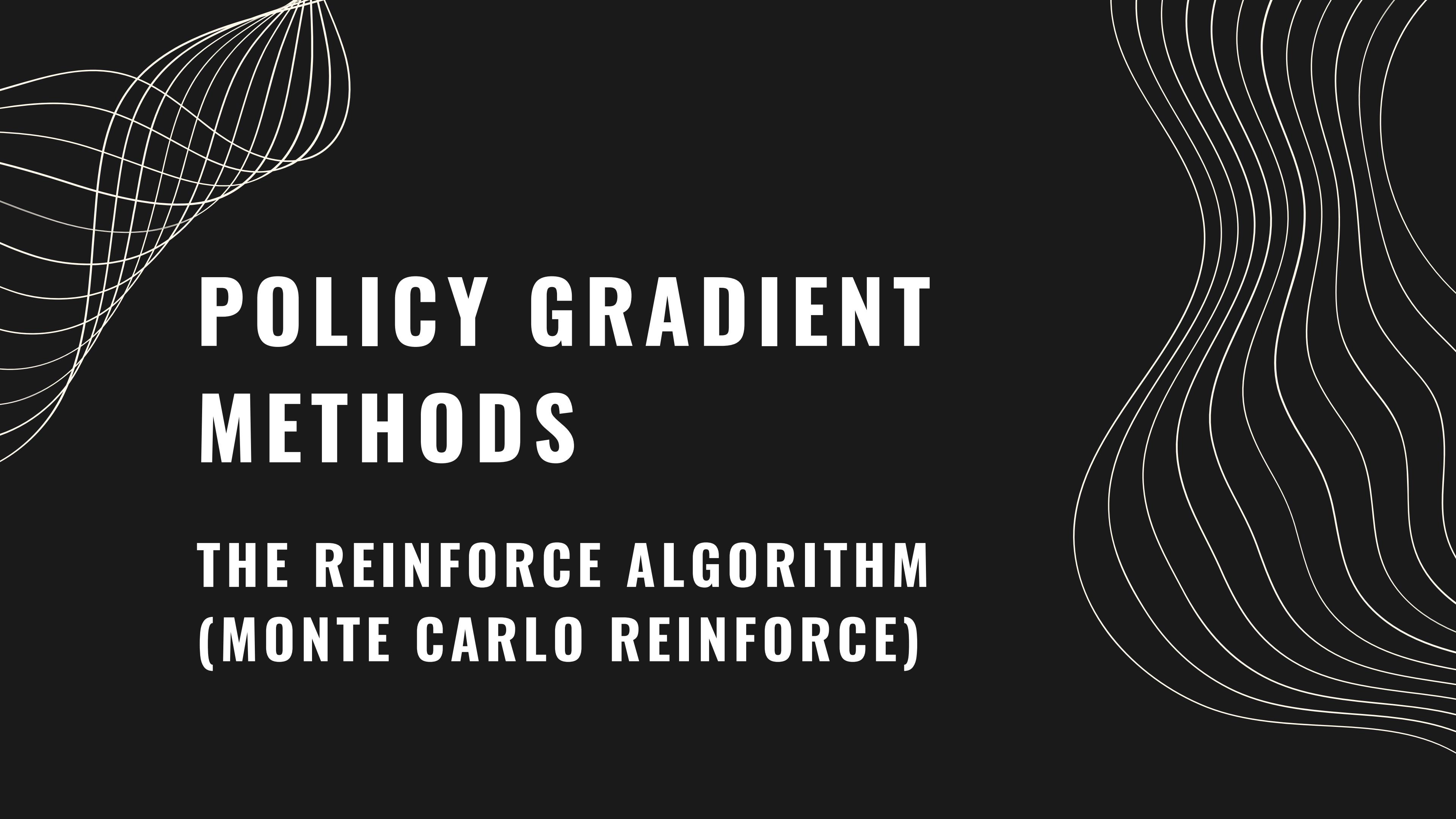
The DQN algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 2
 - Backlog Cost (b): 2
 - Purchase Cost per Unit (p): 10
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The DQN algorithm's policy for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 5
- Backlog Cost (b): 5
- Purchase Cost per Unit (p): 5
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5										
Time from 0 to 4										
-5	-4	-3	-2	-1	0	1	2	3	4	5
0	5	5	5	5	4	3	2	1	0	0
1	5	5	5	5	5	4	3	2	1	0
2	5	5	5	5	5	5	4	3	2	1
3	5	5	5	5	5	4	3	2	1	0
4	5	5	5	5	5	4	3	2	1	0



POLICY GRADIENT METHODS

THE REINFORCE ALGORITHM (MONTE CARLO REINFORCE)

Pseudo code

REINFORCE (Monte-Carlo policy gradient) relies on an estimated return by Monte-Carlo methods using episode samples to update the policy parameter θ . REINFORCE works because the expectation of the sample gradient is equal to the actual gradient:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi}[Q^{\pi}(s, a) \nabla_{\theta} \ln \pi_{\theta}(a|s)] \\ &= \mathbb{E}_{\pi}[G_t \nabla_{\theta} \ln \pi_{\theta}(A_t|S_t)] \quad ; \text{ Because } Q^{\pi}(S_t, A_t) = \mathbb{E}_{\pi}[G_t | S_t, A_t]\end{aligned}$$

Therefore we are able to measure G_t from real sample trajectories and use that to update our policy gradient. It relies on a full trajectory and that's why it is a Monte-Carlo method.

The process is pretty straightforward:

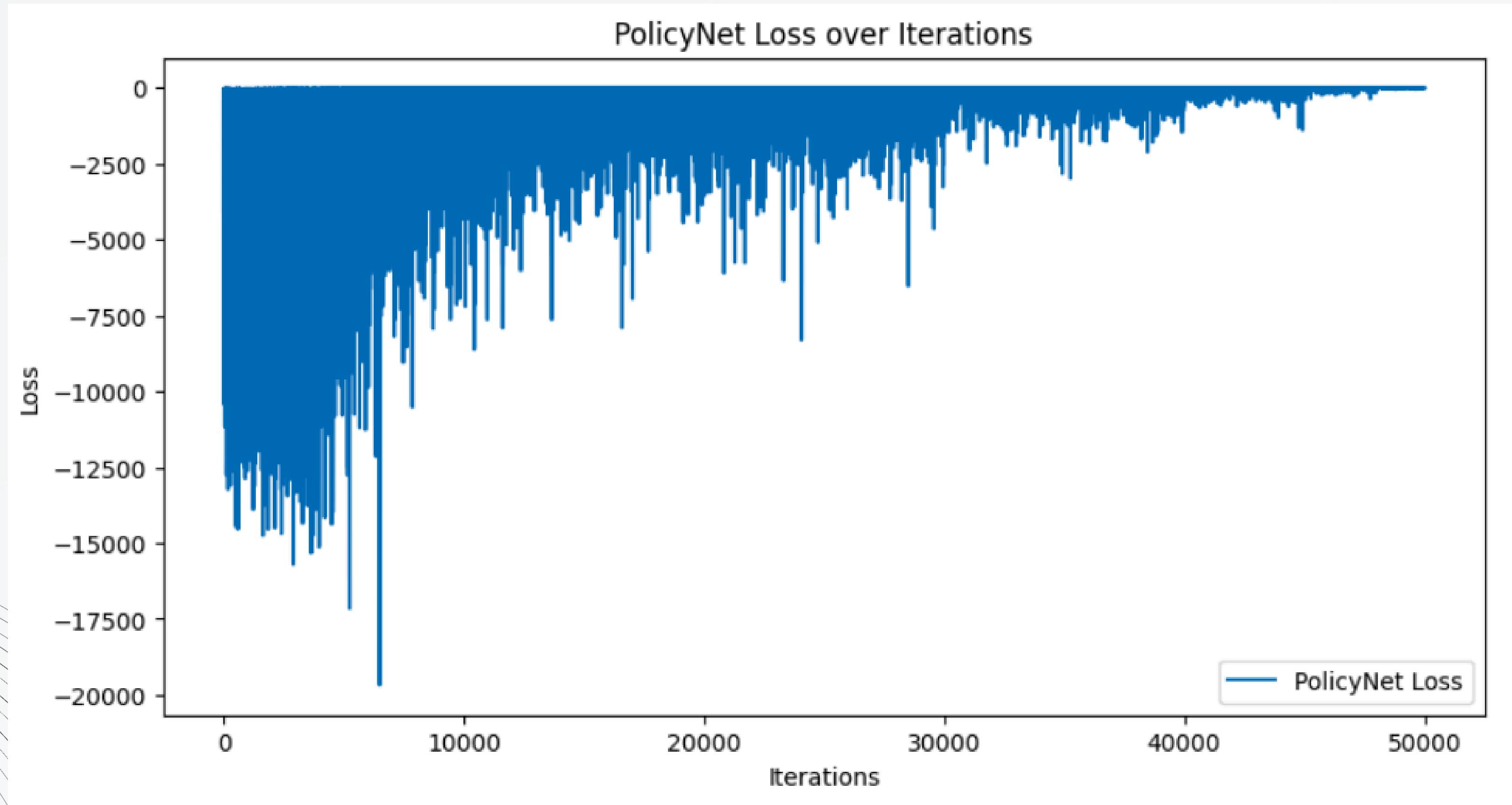
1. Initialize the policy parameter θ at random.
2. Generate one trajectory on policy π_{θ} : $S_1, A_1, R_2, S_2, A_2, \dots, S_T$.
3. For $t=1, 2, \dots, T$:
 1. Estimate the the return G_t ;
 2. Update policy parameters: $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \ln \pi_{\theta}(A_t|S_t)$

The Reinforce algorithm's policy for the given parameters:

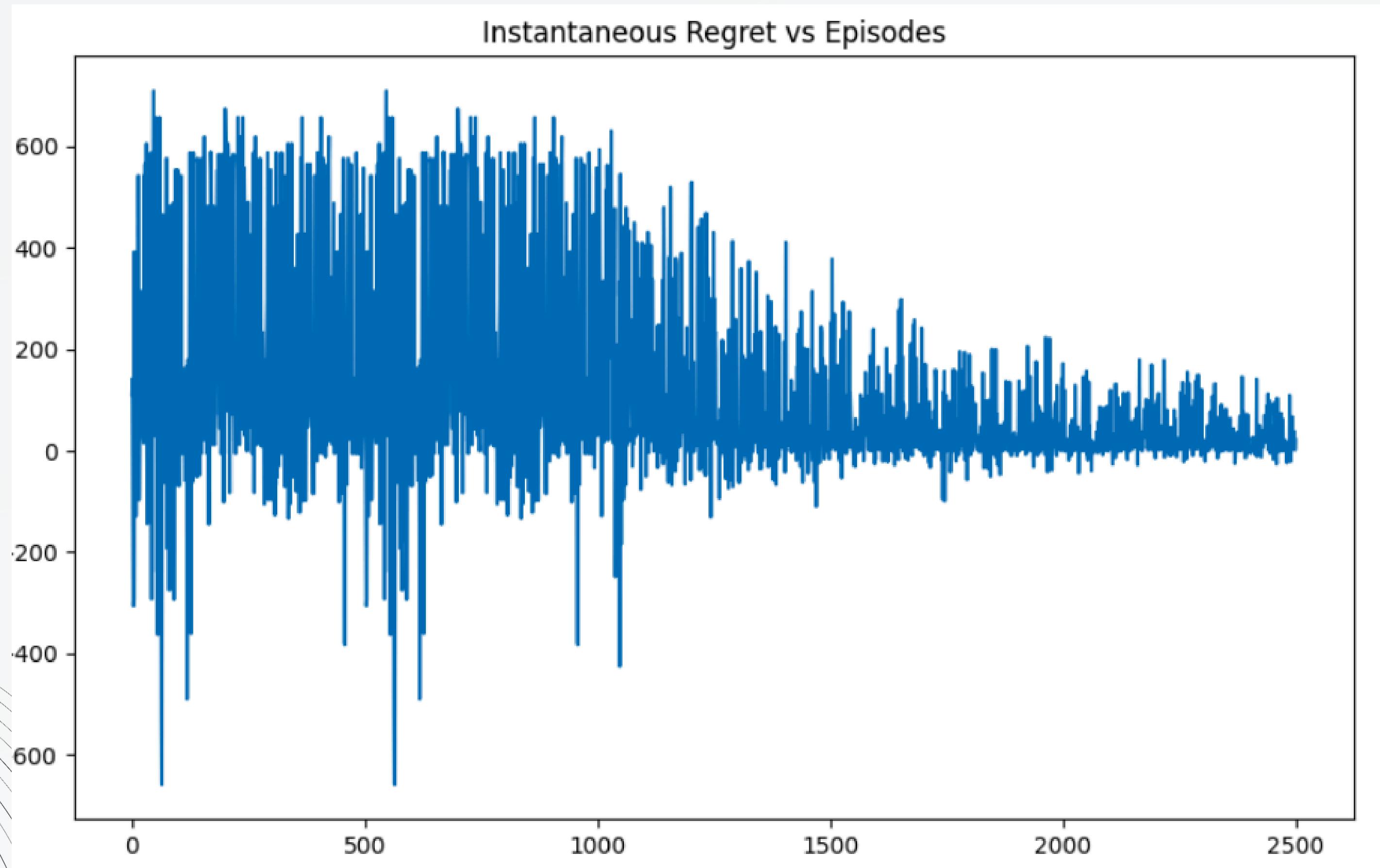
- Horizon (T): 5
- Holding Cost (a): 2
- Backlog Cost (b): 50
- Purchase Cost per Unit (p): 15
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5												
Time from 0 to 4												
-5	-4	-3	-2	-1	0	1	2	3	4	5	4	3
0	5	5	5	5	5	5	5	5	5	5	4	3
1	5	5	5	5	5	5	4	5	5	5	5	5
2	5	5	5	5	5	5	5	5	5	4	3	1
3	5	5	5	5	5	5	5	4	3	2	2	
4	5	5	5	5	4	3	2	2	1	0	0	

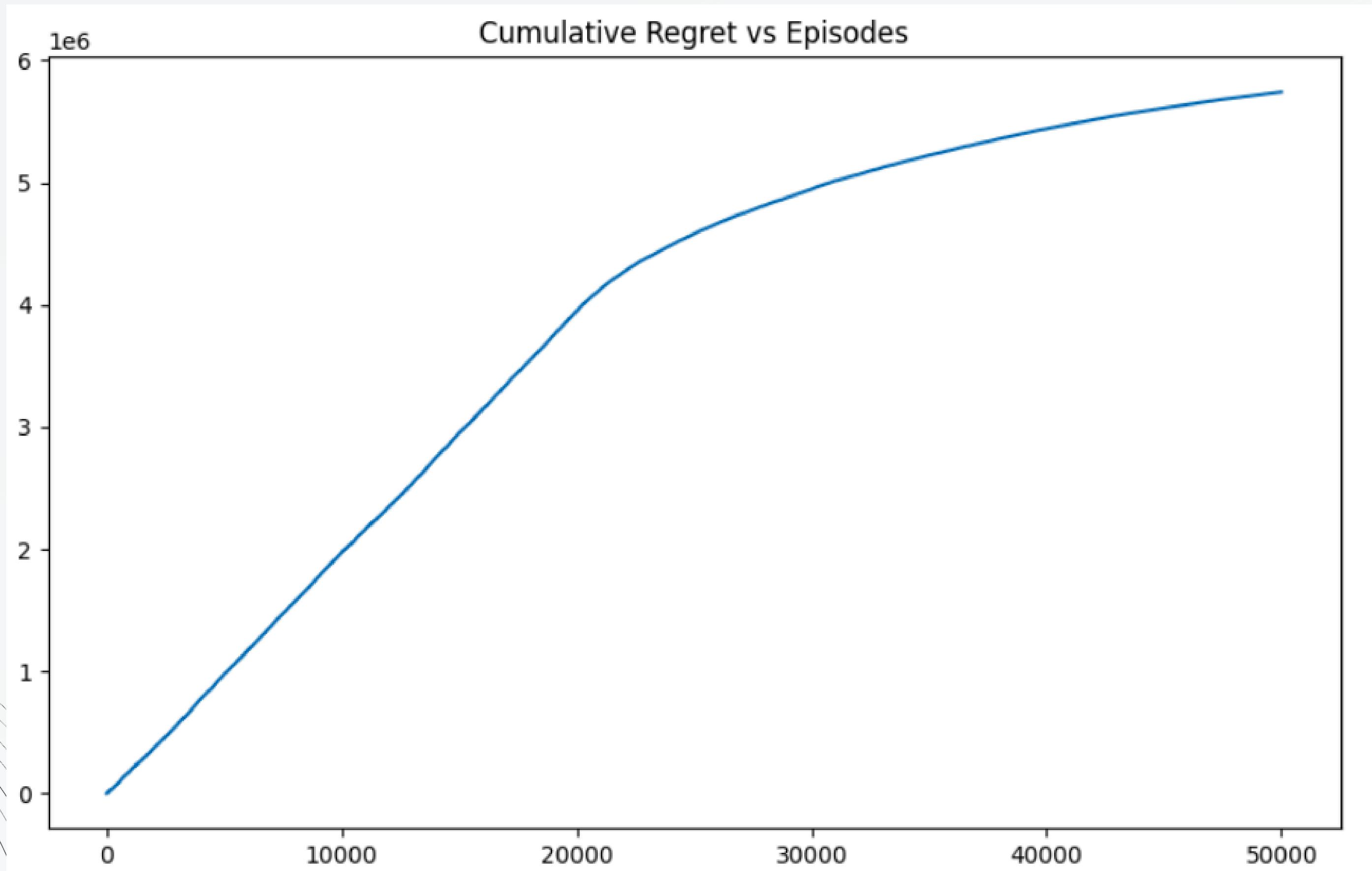
Policy network loss plot

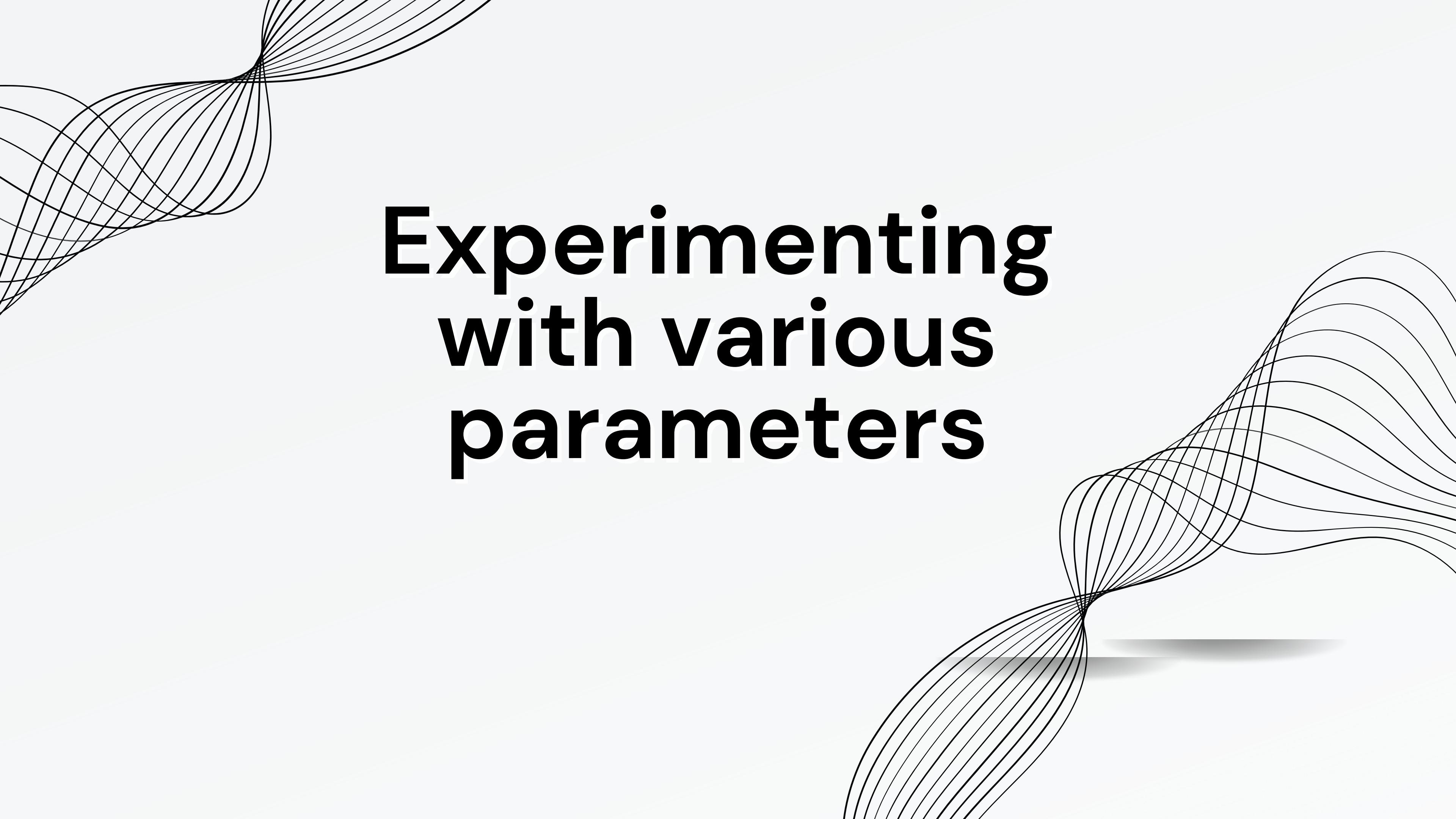


Instantaneous Regret vs Episodes



Cumulative Regret vs Episodes





Experimenting with various parameters

1

The Reinforce algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 3
 - Backlog Cost (b): 3
 - Purchase Cost per Unit (p): 0
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The Reinforce algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 50
 - Backlog Cost (b): 5
 - Purchase Cost per Unit (p): 15
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

States from -5 to 5

	-5	-4	-3	-2	-1	0	1	2	3	4	5	
0	0	4	3	2	1	0	0	0	0	0	0	0
1	1	5	4	3	2	1	0	0	0	0	0	0
2	2	4	3	2	1	0	0	0	0	0	0	0
3	3	5	5	5	5	4	3	2	1	0	0	0
4	4	5	5	5	5	4	3	2	1	0	0	0
5	5	4	3	2	1	0	0	0	0	0	0	0

The Reinforce algorithm's policy for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 2
- Backlog Cost (b): 2
- Purchase Cost per Unit (p): 10
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5											
Time from 0 to 4											
-5	-4	-3	-2	-1	0	1	2	3	4	5	
0	5	5	5	5	4	3	2	1	0	0	
1	5	4	3	2	1	0	0	0	0	0	
2	3	2	1	0	0	0	0	0	0	0	
3	4	3	2	1	0	0	0	0	0	0	
4	4	3	2	1	0	0	0	0	0	0	

The Reinforce algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 5
 - Backlog Cost (b): 5
 - Purchase Cost per Unit (p): 5
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

States from -5 to 5

	-5	-4	-3	-2	-1	0	1	2	3	4	5
0	0	5	5	5	5	5	4	3	2	1	0
1	1	5	5	5	5	5	4	3	2	1	0
2	2	5	5	5	5	5	4	3	2	1	0
3	3	5	4	3	2	1	0	0	0	0	0
4	4	3	2	1	0	0	0	0	0	0	0



**ACTOR
CRITIC**

Pseudo code

Pseudocode of Actor-Critic algorithm[6]

1. Sample { s_t, a_t } using the policy $\pi\theta$ from the actor-network.
2. Evaluate the advantage function A_t . It can be called as TD error δ_t . In Actor-critic algorithm, advantage function is produced by the critic-network.

$$A_{\pi_\theta}(s_t, a_t) = r(s_t, a_t) + V_{\pi_\theta}(s_{t+1}) - V_{\pi_\theta}(s_t)$$

3. Evaluate the gradient using the below expression:

$$\nabla J(\theta) \approx \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t, s_t) A_{\pi_\theta}(s_t, a_t)$$

4. Update the policy parameters, θ

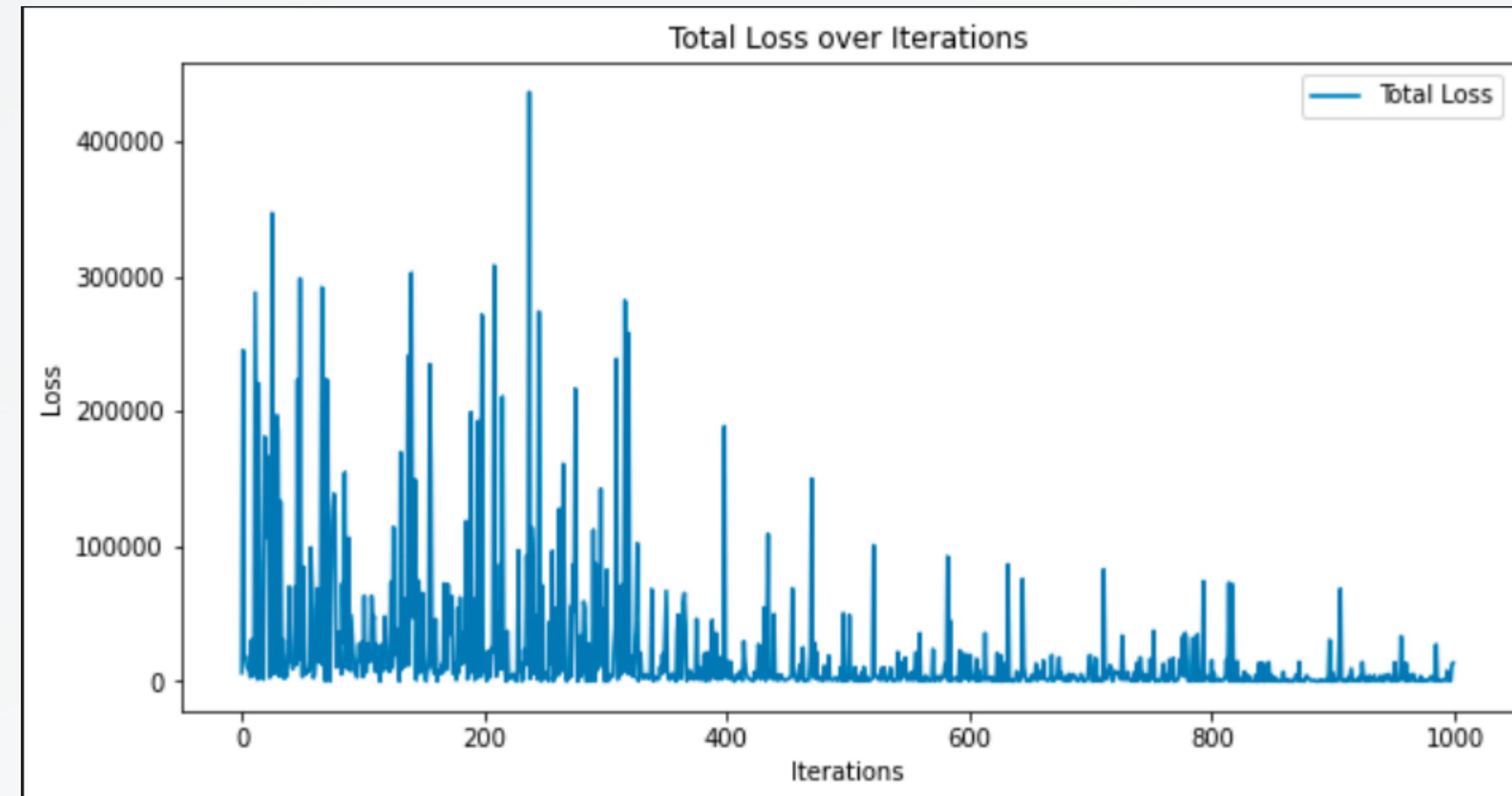
$$\theta = \theta + \alpha \nabla J(\theta)$$

5. Update the weights of the critic based value-based RL(Q-learning). δ_t is equivalent to advantage function.

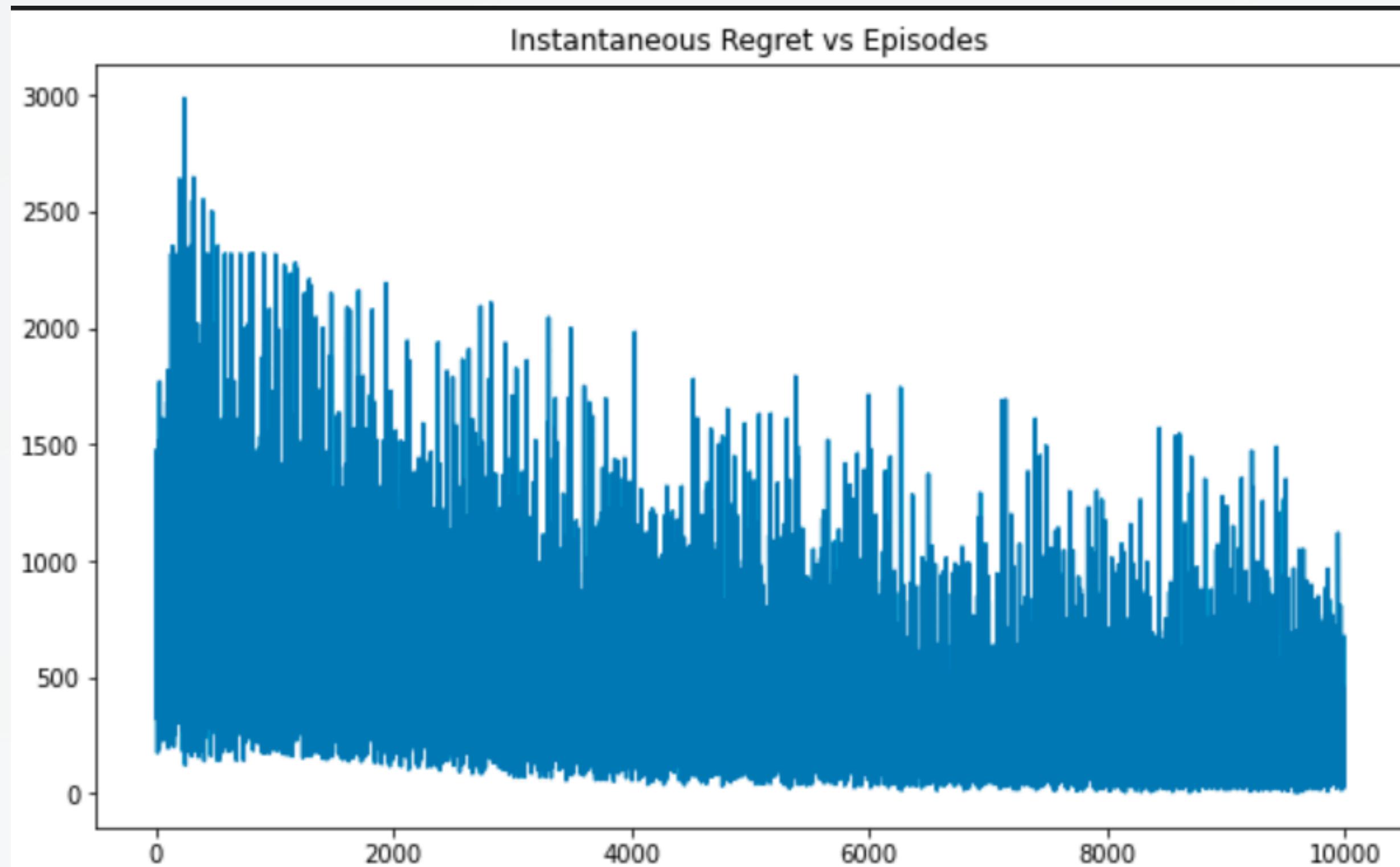
$$w = w + \alpha \delta_t$$

6. Repeat 1 to 5 until we find the optimal policy $\pi\theta$.

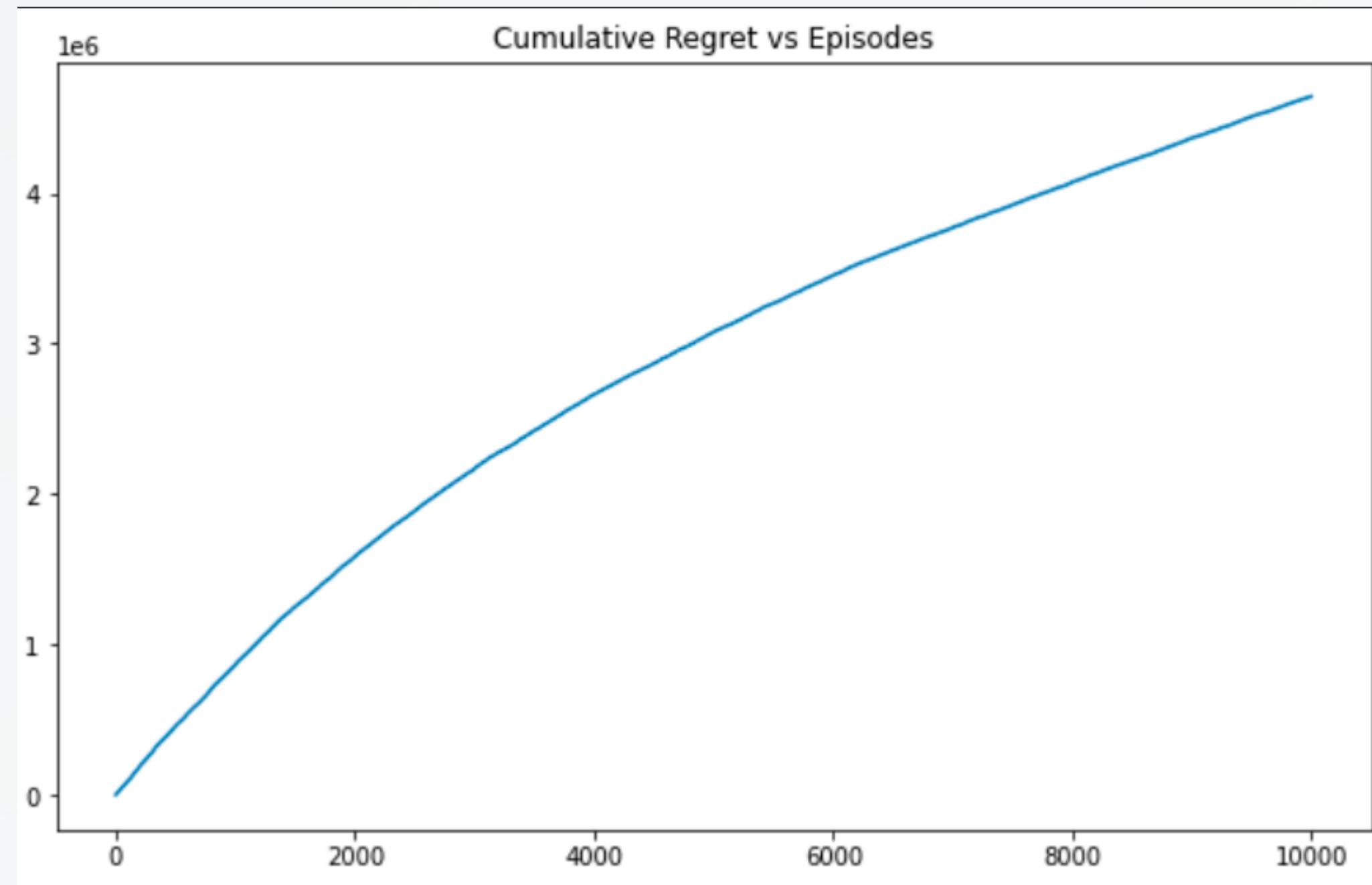
Total Loss Plot

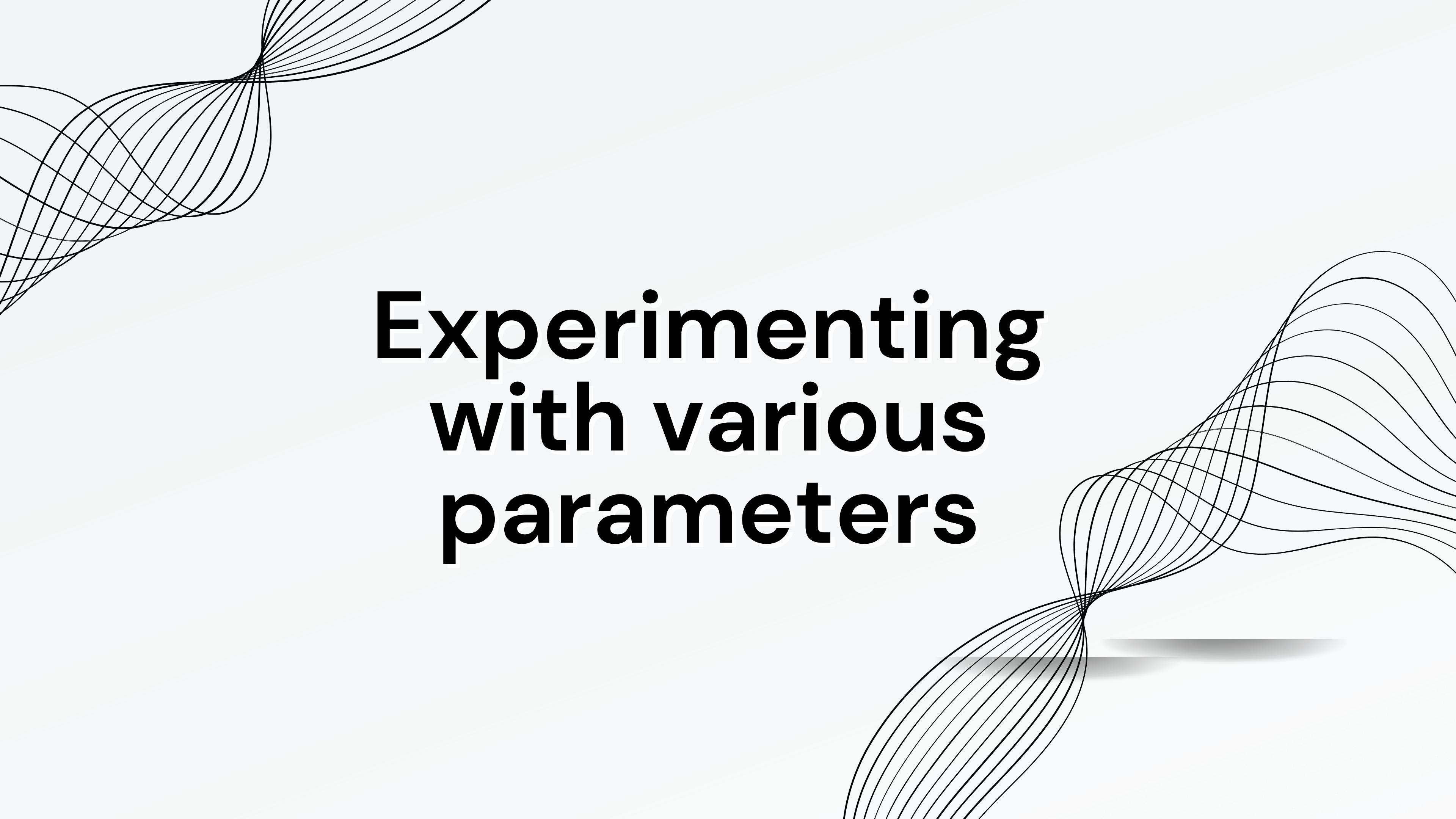


Instantaneous Regret vs Episodes



Cumulative Regret vs Episodes



The background features abstract black line art. On the left, several thin lines radiate from a central point, creating a fan-like or sunburst effect. On the right, a series of concentric, slightly curved lines form a shape reminiscent of a stylized flower or a series of overlapping petals.

Experimenting with various parameters

The Actor Critic algorithm's policy for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 3
- Backlog Cost (b): 3
- Purchase Cost per Unit (p): 0
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5											
Time from 0 to 4											
	-5	-4	-3	-2	-1	0	1	2	3	4	5
0	1	0	0	0	0	0	0	0	0	0	0
1	4	3	2	1	0	0	0	0	0	0	0
2	1	0	0	0	0	0	0	0	0	0	0
3	2	1	0	0	0	0	0	0	0	0	0
4	4	3	2	1	0	0	0	0	0	0	0
5	4	3	2	1	0	0	0	0	0	0	0

The Actor Critic algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 50
 - Backlog Cost (b): 5
 - Purchase Cost per Unit (p): 15
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The Actor Critic algorithm's policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 2
 - Backlog Cost (b): 2
 - Purchase Cost per Unit (p): 10
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The Actor Critic algorithm's policy for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 5
- Backlog Cost (b): 5
- Purchase Cost per Unit (p): 5
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5												
Time from 0 to 4												
-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7
0	5	5	5	5	4	3	2	1	0	0	0	0
1	5	5	5	5	4	3	2	1	0	0	0	0
2	5	5	5	5	5	4	3	2	1	0	0	0
3	5	5	5	5	5	4	3	2	1	0	0	0
4	5	5	5	5	5	4	3	2	1	0	0	0



NATURAL ACTOR CRITIC

**NATURAL GRADIENT WITH FISCHER
PROJECTION MATRIX**

Pseudo code

Input: Parameterized policy $\pi(\mathbf{u}|\mathbf{x}) = p(\mathbf{u}|\mathbf{x}, \theta)$ with initial parameters $\theta = \theta_0$, its derivative $\nabla_\theta \log \pi(\mathbf{u}|\mathbf{x})$.

For $u = 1, 2, 3, \dots$ **do**

For $e = 1, 2, 3, \dots$ **do**

Execute Rollout: Draw initial state $\mathbf{x}_0 \sim p(\mathbf{x}_0)$.

For $t = 1, 2, 3, \dots, N$ **do**

 Draw action $\mathbf{u}_t \sim \pi(\mathbf{u}_t|\mathbf{x}_t)$, observe next state $\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$, and reward $r_t = r(\mathbf{x}_t, \mathbf{u}_t)$.

end.

end.

Critic Evaluation (Episodic): Determine value function

$J = V^\pi(\mathbf{x}_0)$, compatible function approximation $f_w^\pi(\mathbf{x}_t, \mathbf{u}_t)$.

Update: Determine basis functions: $\phi_t = \left[\sum_{t=0}^N \gamma^t \nabla_\theta \log \pi(\mathbf{u}_t|\mathbf{x}_t)^T, 1 \right]^T$;

reward statistics: $R_t = \sum_{t=0}^N \gamma^t r_t$;

Actor-Update: When the natural gradient is converged,

$\angle(\mathbf{w}_{t+1}, \mathbf{w}_{t-\tau}) \leq \epsilon$, update the policy parameters: $\theta_{t+1} = \theta_t + \alpha \mathbf{w}_{t+1}$.

6: **end.**

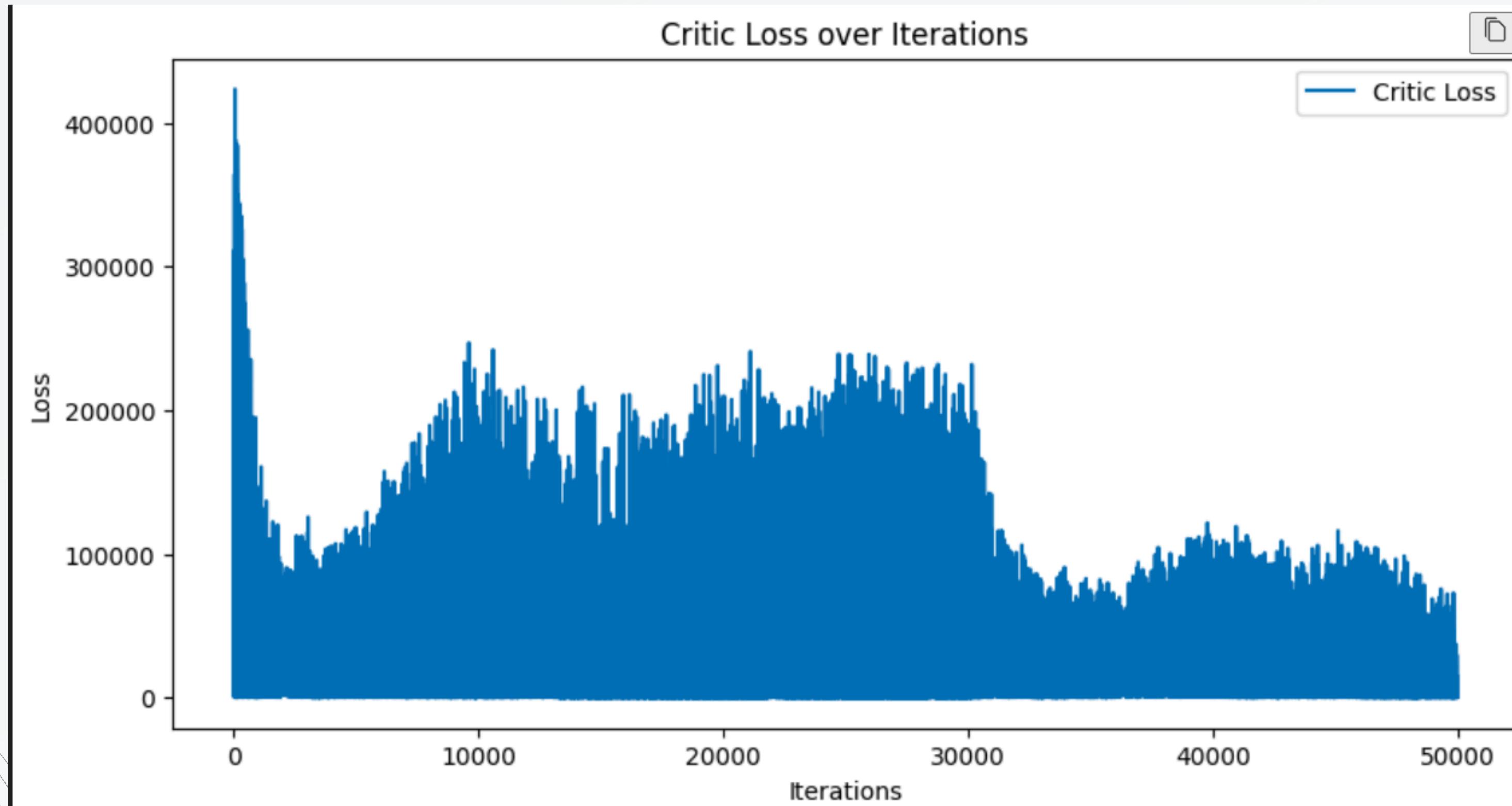
The Natural Actor critic policy for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 2
- Backlog Cost (b): 50
- Purchase Cost per Unit (p): 15
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

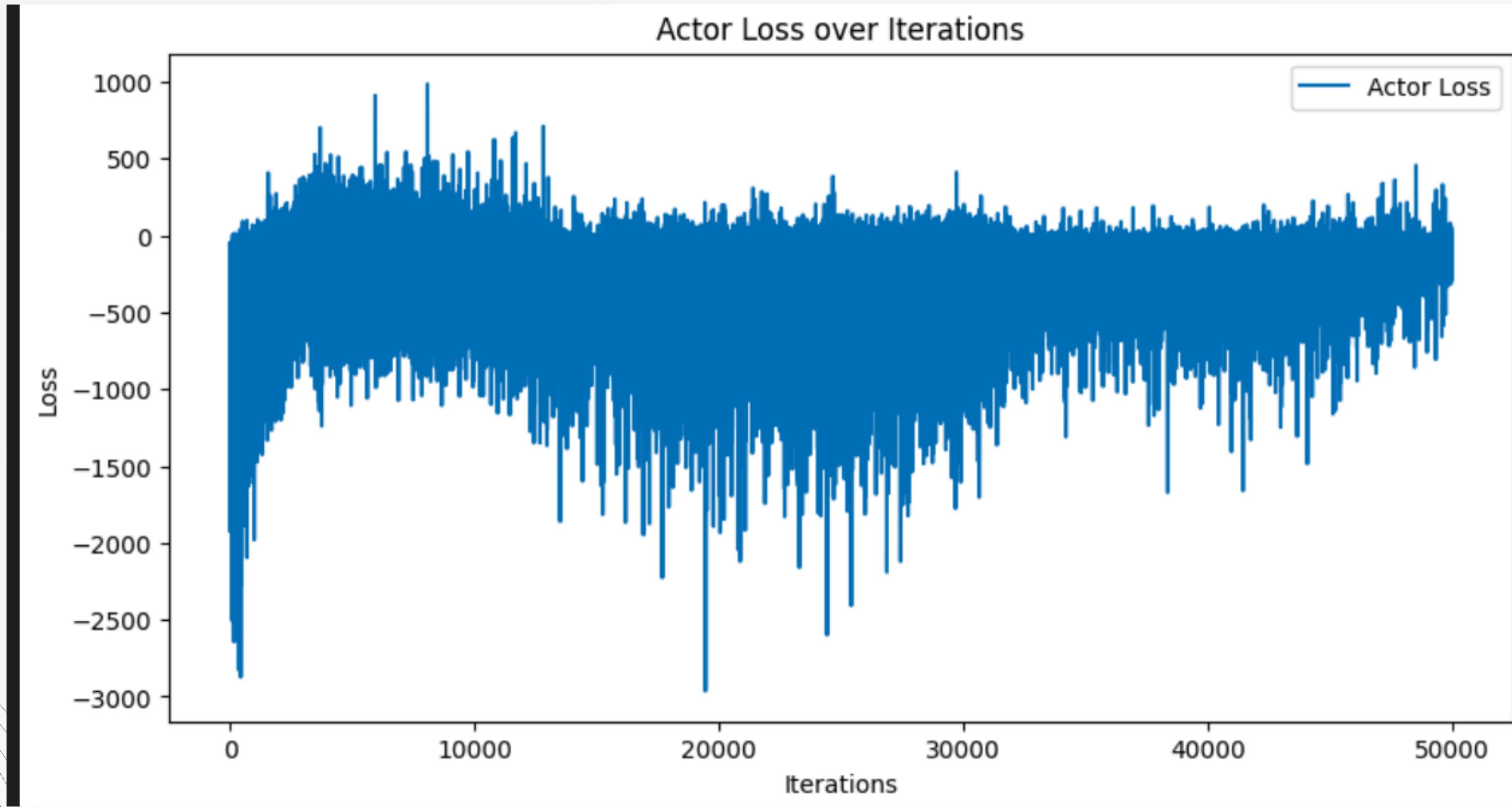
States from -5 to 5

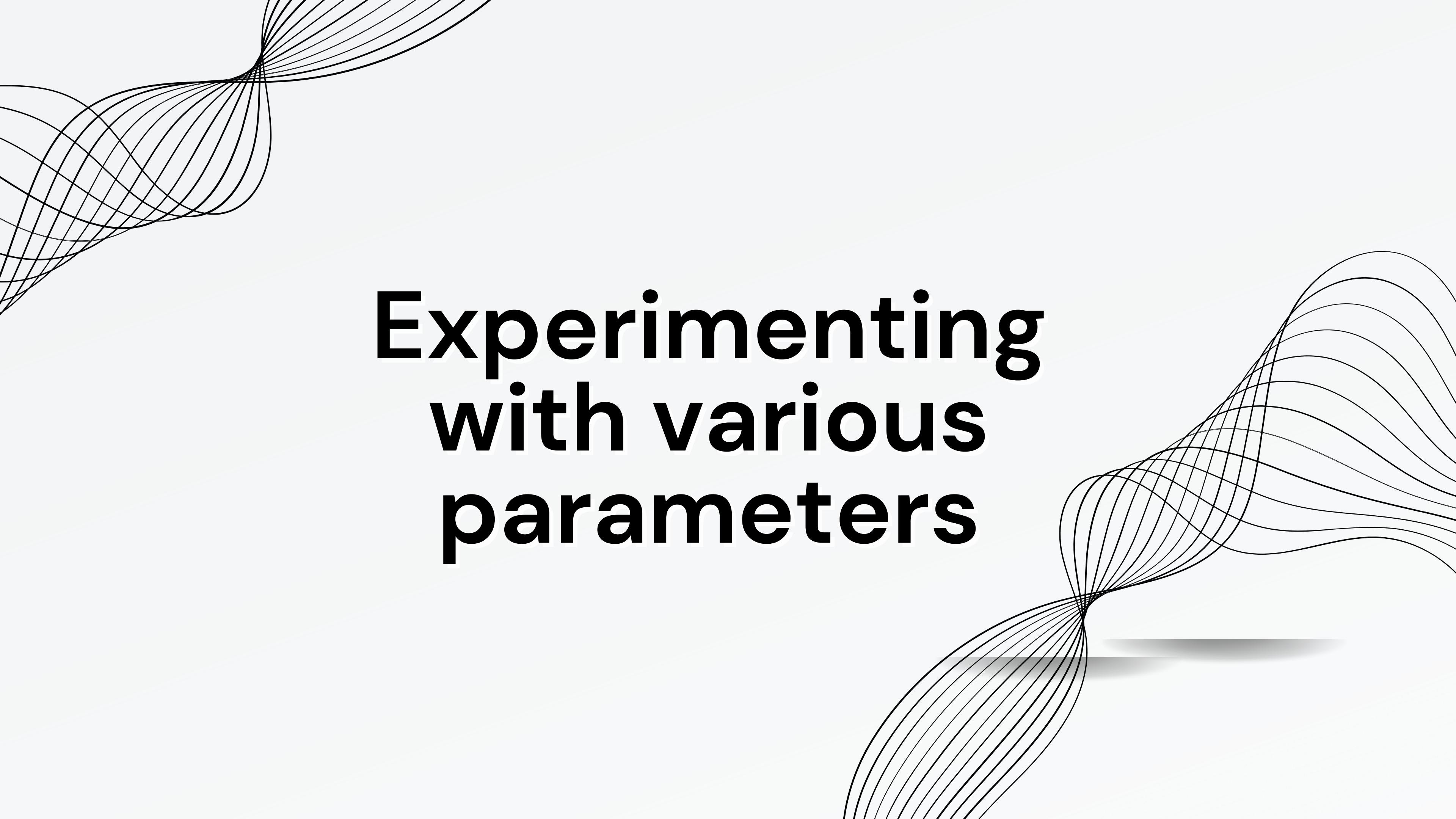
Time from 0 to 5													
		-5	-4	-3	-2	-1	0	1	2	3	4	5	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+													
	0	5	5	5	5	5	5	5	5	5	4		
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+		1	5	5	5	5	5	5	4	3	2	1	0
	2	4	5	5	4	5	5	5	5	5	4	3	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+		3	5	5	5	5	5	5	4	3	2	1	
	4	5	5	5	5	4	3	3	1	0	0	0	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+		5	5	5	4	3	2	1	0	0	0	0	

Critic loss plot



Actor network loss plot



The background features abstract black line art. On the left, several thin lines radiate from a central point, creating a fan-like or sunburst effect. On the right, a series of concentric, slightly curved lines form a shape reminiscent of a stylized flower or a series of overlapping petals.

Experimenting with various parameters

1

The Natural Actor-Critic policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 3
 - Backlog Cost (b): 3
 - Purchase Cost per Unit (p): 0
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The Natural Actor-Critic policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 50
 - Backlog Cost (b): 5
 - Purchase Cost per Unit (p): 15
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

The Natural Actor Critic policy for the given parameters:

- Horizon (T): 5
 - Holding Cost (a): 2
 - Backlog Cost (b): 2
 - Purchase Cost per Unit (p): 10
 - State Space (S): -5 to 5
 - Action Space: 0 to 5
 - Demand Range: 0 to 5

States from -5 to 5

The Natural Actor-Critic algorithm's policy
for the given parameters:

- Horizon (T): 5
- Holding Cost (a): 5
- Backlog Cost (b): 5
- Purchase Cost per Unit (p): 5
- State Space (S): -5 to 5
- Action Space: 0 to 5
- Demand Range: 0 to 5

States from -5 to 5												
Time from 0 to 5												
	-5	-4	-3	-2	-1	0	1	2	3	4	5	
0	5	5	5	5	4	3	2	1	0	0	0	
1	5	4	3	2	1	0	0	0	0	0	0	
2	5	5	5	5	4	3	2	1	0	0	0	
3	5	5	5	5	4	3	2	1	0	0	0	
4	5	5	5	4	3	2	1	0	0	0	0	
5	5	5	4	3	2	1	0	0	0	0	0	

THANK YOU

