



TEXAS TECH UNIVERSITY SYSTEM™



Reinforcement Learning Algorithms

Code Implementation

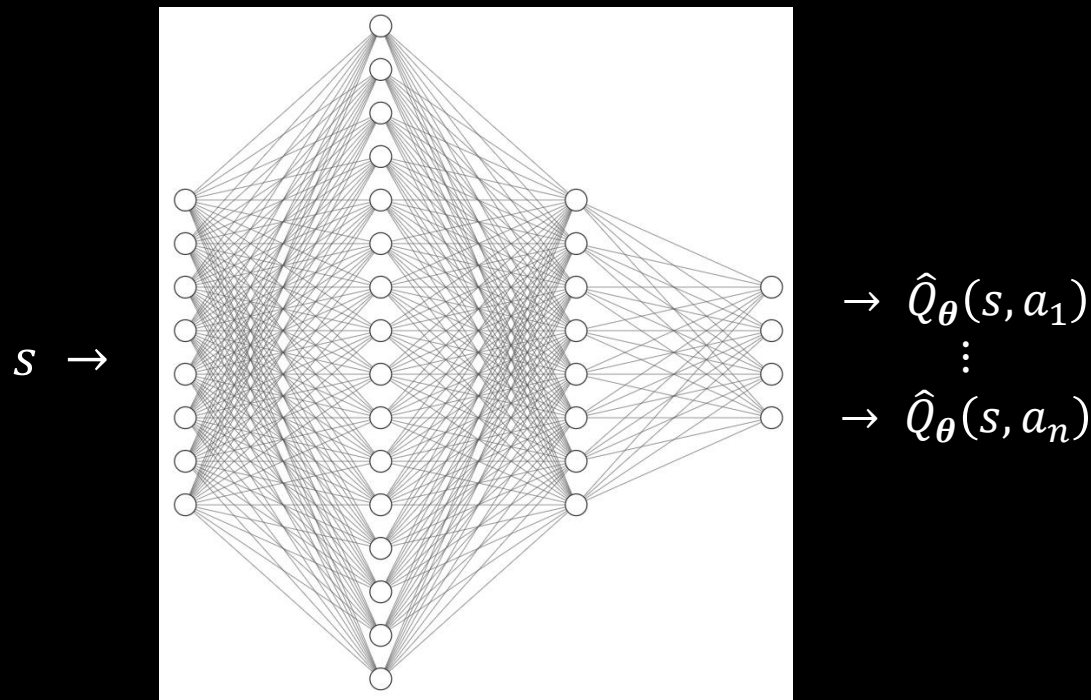
Saif Al Wahaibi
PhD Candidate

July 19th, 2023

Q-Learning

- Intuition:

- Estimate $Q_{\pi}(s, a)$ via function approximation



- Math:

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi} \left[\left(\delta_{TD} - \hat{Q}_{\theta}(s, a) \right)^2 \right] = \mathbb{E}_{\pi} \left[\left(r + \gamma \max_{a'} \hat{Q}_{\theta}(s', a') - \hat{Q}_{\theta}(s, a) \right)^2 \right]$$



Q-Learning Algorithm

- Pseudocode:

Initialize $\hat{Q}_\theta(s, a)$ with random weights

for episode = 1, 2, 3, ..., E do

Initialize environment s_0

for $t = 0, 1, 2, \dots, T$ do

Select action a_t randomly with probability ϵ , otherwise

$$a_t = \underset{a_t}{\operatorname{argmax}} \hat{Q}_\theta(s_t, a_t)$$

*Execute action a_t in environment and observe r_{t+1} , s_{t+1} ,
and terminal or truncate flags*

Set TD target $\delta_{TD} = r_{t+1}$ if terminal or truncate, otherwise

$$\delta_{TD} = r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}_\theta(s_{t+1}, a_{t+1})$$

Perform a gradient descent step on

$$J(\theta) = \mathbb{E}_\pi \left[\left(\delta_{TD} - \hat{Q}_\theta(s_t, a_t) \right)^2 \right]$$

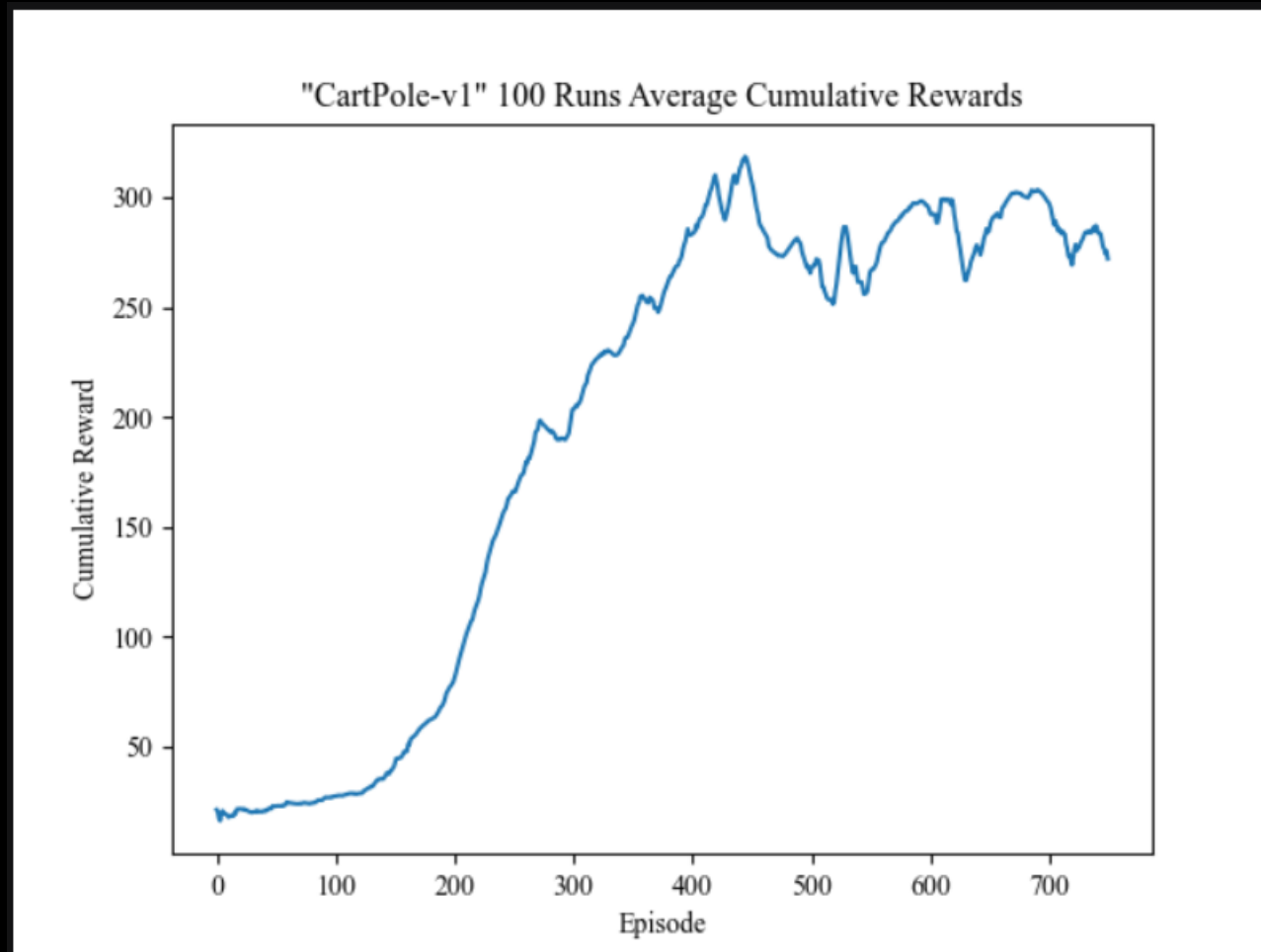
Set s_{t+1} as current state

end for

end for

Q-Learning Example

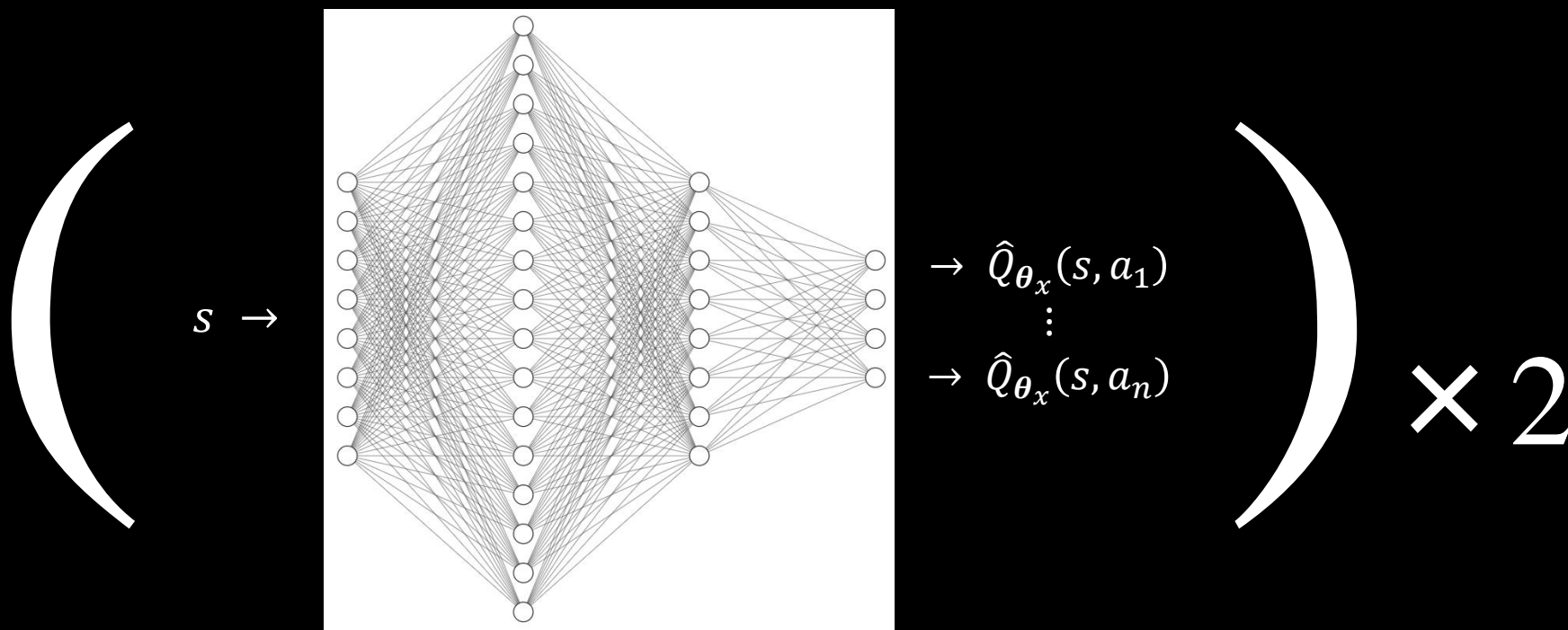
- “CartPole-v1” Gym Environment:



Double Q-Learning

- Intuition:

- Estimate $Q_{\pi}(s, a)$ via function approximation
- Use two Q-networks to handle maximization bias





Double Q-Learning

- Example of Maximization Bias:

- Assume $\hat{Q} = 150$
- 10 noisy estimates of \hat{Q} :

1	2	3	4	5	6	7	8	9	10
150	151	151	150	151	151	151	152	150	149

- Another 10 noisy estimates of \hat{Q} generated similarly but independently:

1	2	3	4	5	6	7	8	9	10
150	151	147	151	151	150	150	150	151	148

- Math:

$$\mathcal{J}(\theta_i) = \mathbb{E}_{\pi} \left[\left(r + \gamma \hat{Q}_{\theta_j} \left(s', \underset{a'}{\operatorname{argmax}} \hat{Q}_{\theta_i}(s', a') \right) - \hat{Q}_{\theta_i}(s, a) \right)^2 \right]$$



Double Q-Learning Algorithm

- Pseudocode:

Initialize both $\hat{Q}_{\theta_x}(s, a)$ with random weights

for episode = 1, 2, 3, ..., E do

Initialize environment s_0

for $t = 0, 1, 2, \dots, T$ do

Select action a_t randomly with probability ϵ , otherwise

$$a_t = \underset{a_t}{\operatorname{argmax}} \left(\frac{\hat{Q}_{\theta_1}(s_t, a_t) + \hat{Q}_{\theta_2}(s_t, a_t)}{2} \right)$$

*Execute action a_t in environment and observe r_{t+1}, s_{t+1} ,
and terminal or truncate flags*

Choose at random either to update 1 or 2

if $i = 1$ or 2 then

*Set TD target $\delta_{TD} = r_{t+1}$ if terminal or truncate,
otherwise*

$$a_{t+1}^* = \underset{a_{t+1}}{\operatorname{argmax}} \hat{Q}_{\theta_i}(s_{t+1}, a_{t+1})$$

$$\delta_{TD} = r_{t+1} + \gamma \hat{Q}_{\theta_{3-i}}(s_{t+1}, a_{t+1}^*)$$



Double Q -Learning Algorithm

$$a^*_{t+1} = \underset{a_{t+1}}{\operatorname{argmax}} \hat{Q}_{\theta_i}(s_{t+1}, a_{t+1})$$

$$\delta_{TD} = r_{t+1} + \gamma \hat{Q}_{\theta_{3-i}}(s_{t+1}, a^*_{t+1})$$

Perform a gradient descent step on

$$J(\theta) = \mathbb{E}_{\pi} \left[\left(\delta_{TD} - \hat{Q}_{\theta_i}(s_t, a_t) \right)^2 \right]$$

end if

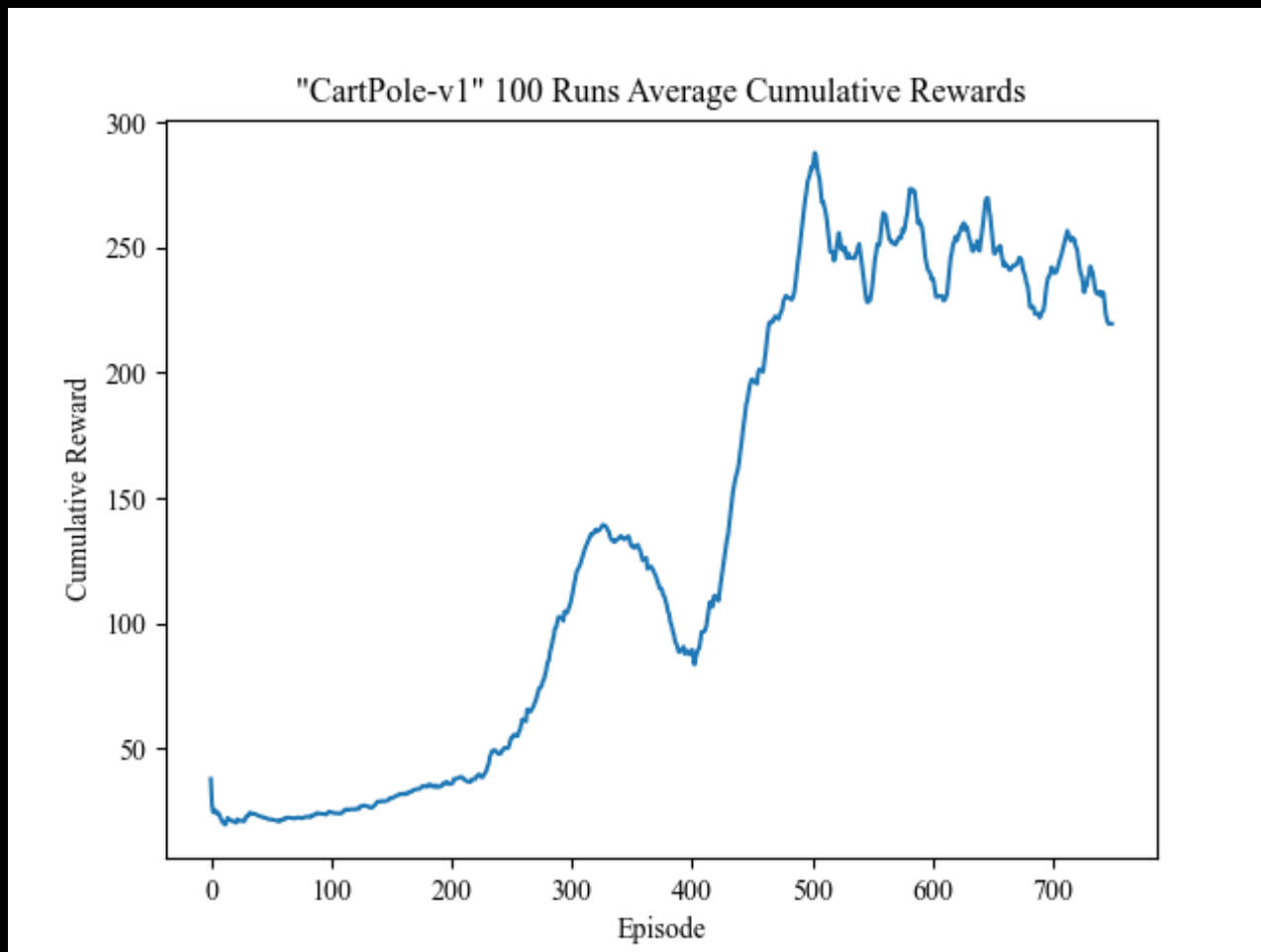
Set s_{t+1} as current state

end for

end for

Double Q -Learning Example

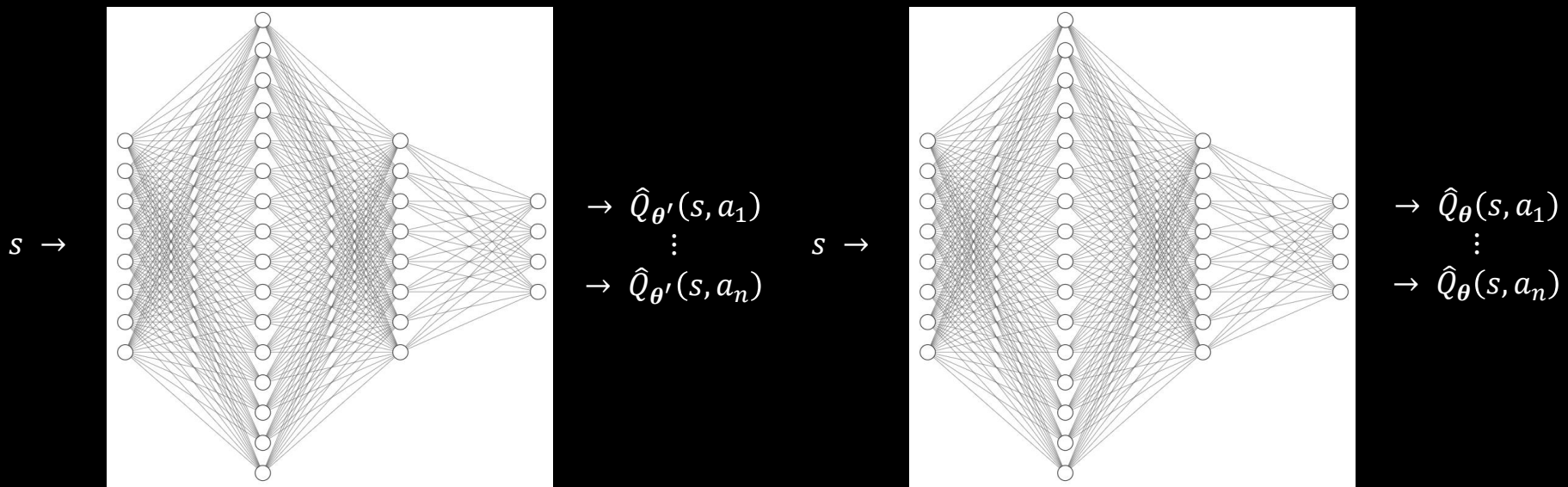
- “CartPole-v1” Gym Environment:



Deep Q-Learning Network

Intuition:

- Estimate $Q_{\pi}(s, a)$ via function approximation
- Use two Q-networks, one frozen and the other online, to fix the TD-target
- Utilize a replay buffer to store experience and train networks in randomized batches



Math:

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi} \left[\left(r + \gamma \max_{a'} \hat{Q}_{\theta'}(s', a') - \hat{Q}_{\theta}(s, a) \right)^2 \right]$$



Deep Q-Learning Network Algorithm

- Pseudocode:

Initialize both $\hat{Q}_\theta(s, a)$ & $\hat{Q}_{\theta'}(s, a)$ with random weights

Initialize replay buffer

for episode = 1, 2, 3, ..., E do

Initialize environment s_0

for $t = 0, 1, 2, \dots, T$ do

Select action a_t randomly with probability ϵ , otherwise

$$a_t = \underset{a_t}{\operatorname{argmax}} \hat{Q}_\theta(s_t, a_t)$$

*Execute action a_t in environment and observe r_{t+1} , s_{t+1} ,
 and terminal or truncate flags*

*Save s_t , a_t , r_{t+1} , s_{t+1} , and terminal or truncate flags in
 replay buffer*

*Sample a random batch of $s_t^{\mathcal{R}}$, $a_t^{\mathcal{R}}$, $r_{t+1}^{\mathcal{R}}$, $s_{t+1}^{\mathcal{R}}$, and
 terminal or truncate flags from replay buffer*

Set TD target $\delta_{TD} = r_{t+1}^{\mathcal{R}}$ if terminal or truncate, otherwise

$$\delta_{TD} = r_{t+1}^{\mathcal{R}} + \gamma \max_{a^{\mathcal{R}}_{t+1}} \hat{Q}_{\theta'}(s_{t+1}^{\mathcal{R}}, a^{\mathcal{R}}_{t+1})$$



Deep Q-Learning Network Algorithm

Set TD target $\delta_{TD} = r^{\mathcal{R}}_{t+1}$ if terminal or truncate, otherwise
$$\delta_{TD} = r^{\mathcal{R}}_{t+1} + \gamma \max_{a^{\mathcal{R}}_{t+1}} \hat{Q}_{\theta'}(s^{\mathcal{R}}_{t+1}, a^{\mathcal{R}}_{t+1})$$

Perform a gradient descent step on

$$J(\theta) = \mathbb{E}_{\pi} \left[\left(\delta_{TD} - \hat{Q}_{\theta}(s^{\mathcal{R}}_t, a^{\mathcal{R}}_t) \right)^2 \right]$$

Count learning steps:

$$i = i + 1$$

if $i = \text{replace}$ then

Update target network:

$$\theta' \leftarrow \theta$$

end if

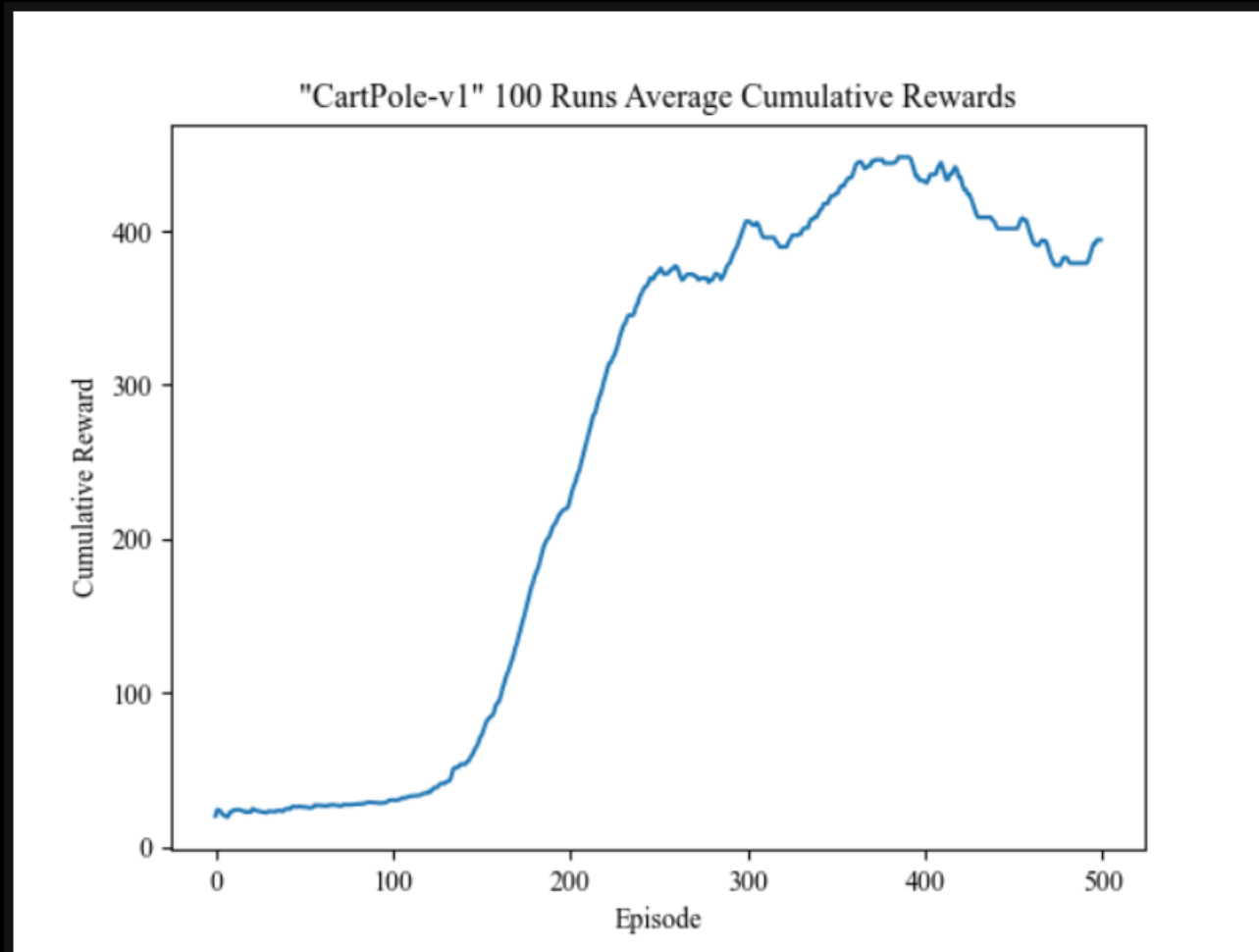
Set s_{t+1} as current state

end for

end for

Deep Q -Learning Network Example

- “CartPole-v1” Gym Environment:





Double Deep Q-Learning Network

- Intuition:

- Estimate $Q_{\pi}(s, a)$ via function approximation
- Use two Q-networks, one frozen and the other online, to fix the TD-target
- Utilize a replay buffer to store experience and train networks in randomized batches
- Use the online and target Q-networks to handle maximization bias

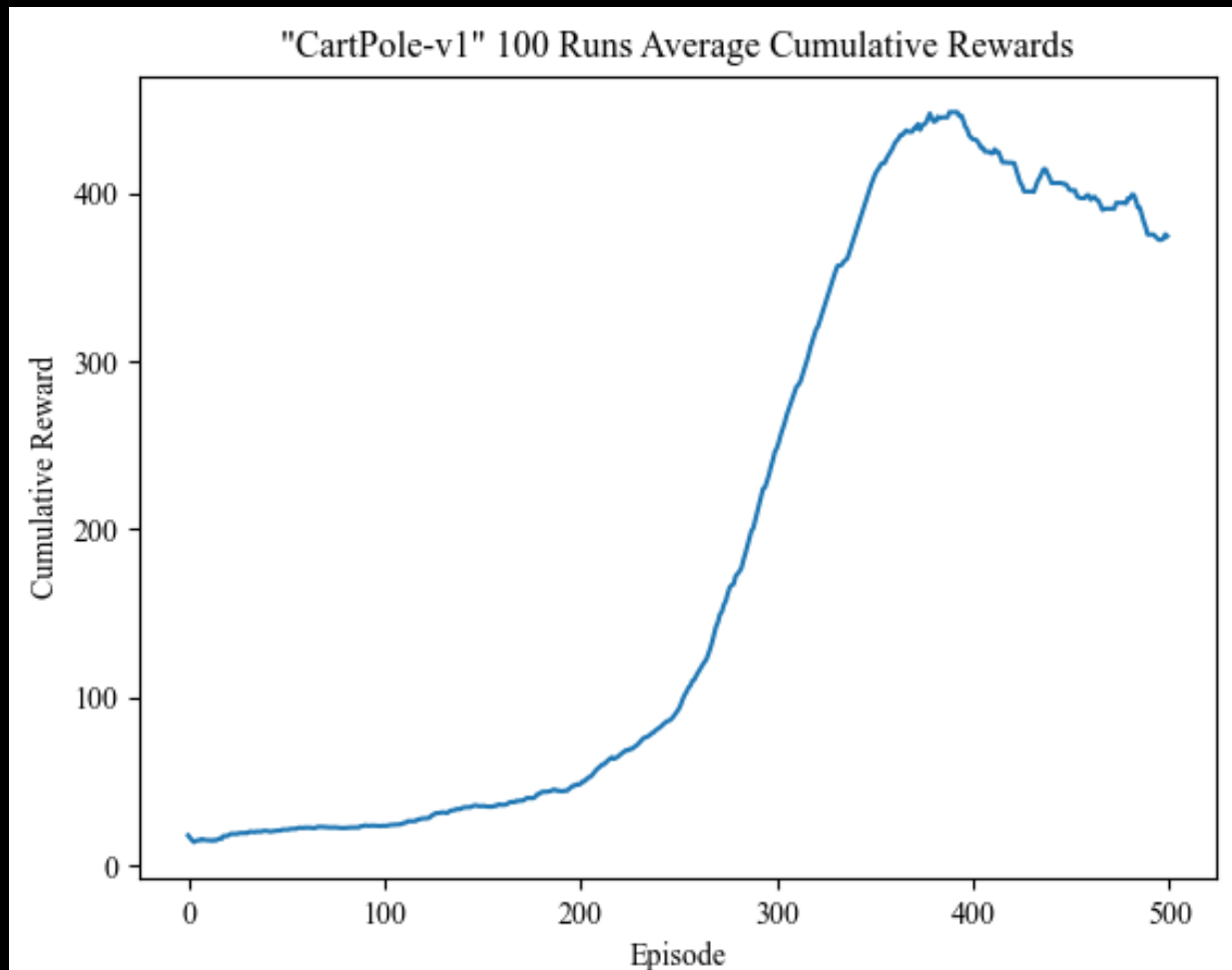
- Math:

$$\mathcal{J}(\boldsymbol{\theta}) = \mathbb{E}_{\pi} \left[\left(r + \gamma \hat{Q}_{\boldsymbol{\theta}'} \left(s', \underset{a'}{\operatorname{argmax}} \hat{Q}_{\boldsymbol{\theta}}(s', a') \right) - \hat{Q}_{\boldsymbol{\theta}}(s, a) \right)^2 \right]$$



Double Deep Q -Learning Network Example

- “CartPole-v1” Gym Environment:





TEXAS TECH UNIVERSITY SYSTEM™