

# Reinforcement Learning

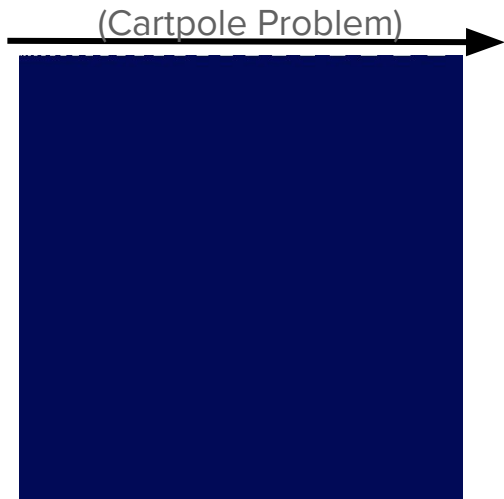
---

Alexis W. Mills, Daniel Chai, Weishi Yan

# Reinforcement Learning Overview

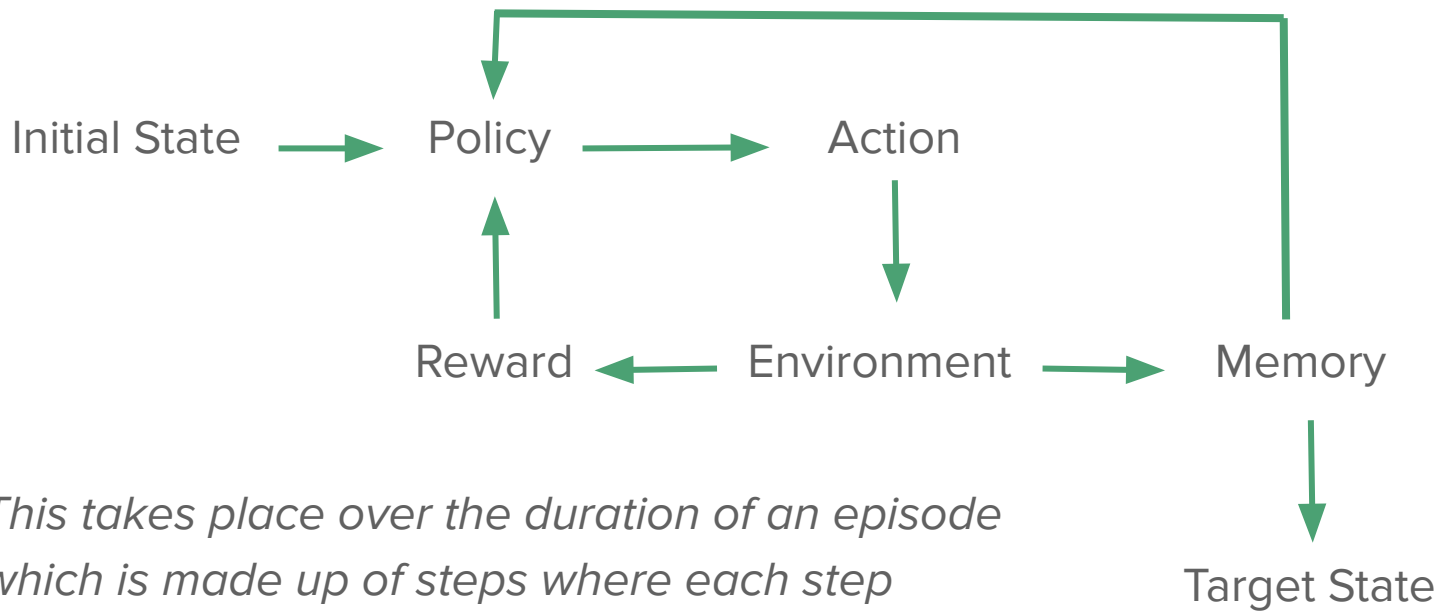
The system is made up of 5 key components:

1. Agent
2. Environment
3. Action
4. Reward
5. State



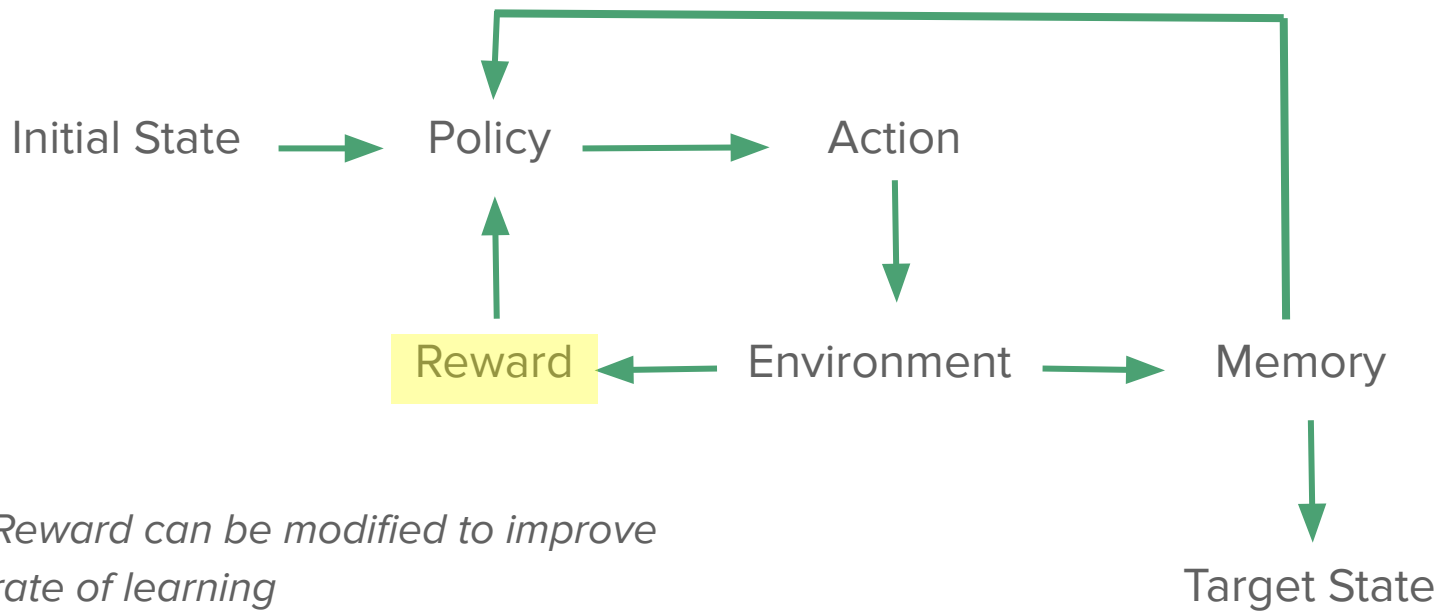
1. Controls the cart (black square)
2. Angle of pole and position of cart
3. Move in positive or negative direction along 1D surface
4. Reward=1 if agent has not failed (moved cart beyond threshold or allow pole angle to fall beyond threshold)
5. Consists of cart position, cart velocity, pole angle, and pole velocity at tip

# Learning Algorithm



*This takes place over the duration of an episode which is made up of steps where each step allows the agent to take an action*

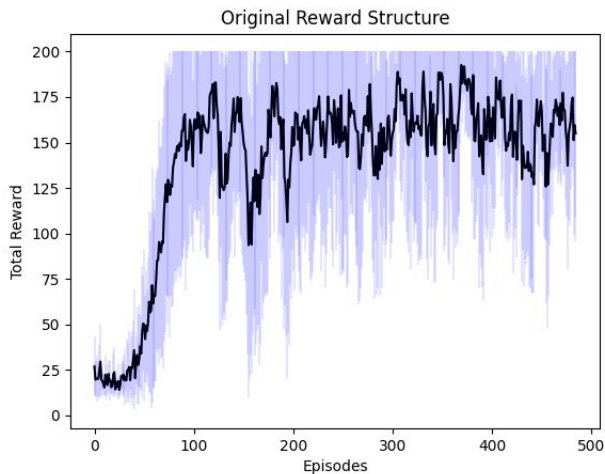
# Learning Algorithm



# Reward Structure

## *Original*

A reward of value 1.0 is granted so long as the pole angle is  $-41.8^\circ < \theta < 41.8^\circ$  and the cart position is  $-2.4 < x < 2.4$



## *Modified*

The reward value is varied depending on how close the agent comes to failing during the step

*Calculating the Reward:*

$$\theta_r = (\theta_t - \theta_a) / \theta_t$$

$\theta_r$ :  $\theta$  reward

$\theta_t$ :  $41.8^\circ$

$\theta_a$ : Absolute value of  $\theta$   
angle (from current state)

$$x_r = (x_t - x_a) / x_t$$

$x_r$ : x reward

$x_t$ : 2.4

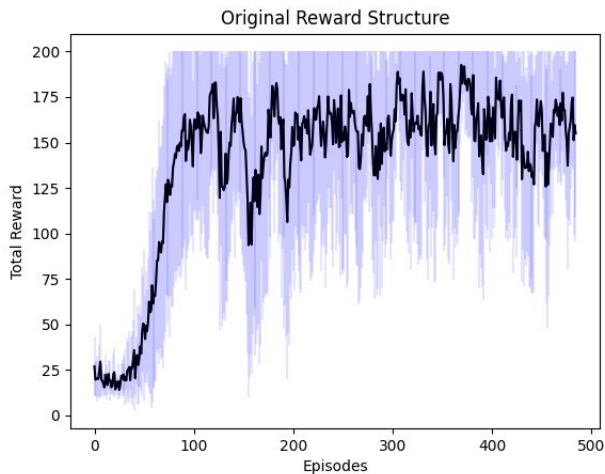
$x_a$ : Absolute value of x  
position (from current state)

$$\text{Modified reward} = (\theta_r + x_r) / 2$$

# Reward Structure

## *Original*

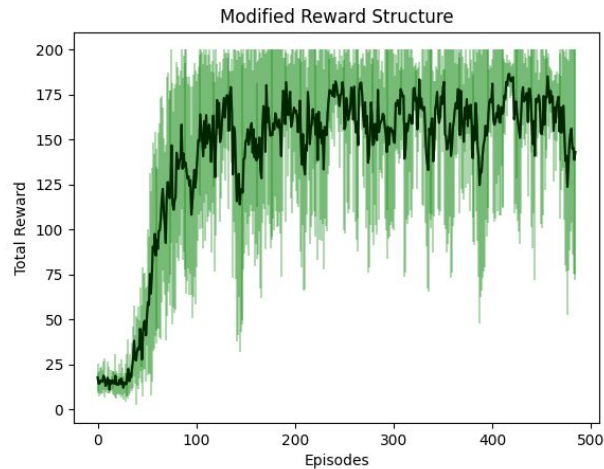
A reward of value 1.0 is granted so long as the pole angle is  $-41.8^\circ < \theta < 41.8^\circ$  and the cart position is  $-2.4 < x < 2.4$



## *Modified*

The reward value is varied depending on how close the agent comes to failing during the step

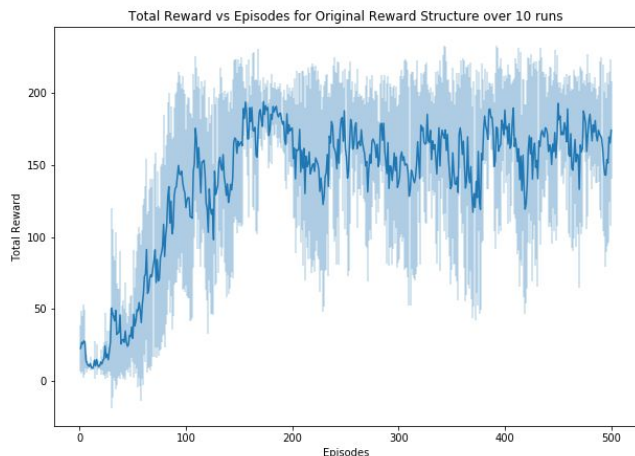
$$\text{Modified reward} = (\theta_r + x_r)/2$$



# Reward Structure

## *Original*

A reward of value 1.0 is granted so long as the pole angle is  $-41.8^\circ < \theta < 41.8^\circ$  and the cart position is  $-2.4 < x < 2.4$



## *Modified\_version2*

Since the goal of the problem is to have the pole stay upright as long as possible, physics equation was used to calculate the time the pole would stay up.

Using this physics equation:

$$x = x_{\text{initial}} + v * t$$

Transformed into:

$$\Theta_{\text{time}} = (\Theta_{\text{threshold}} - \text{abs}(\Theta)) / \Theta_{\text{dot}}$$

Reward Calculation - hyperbolic tangent (tanh):

$$\Theta_{\text{reward}} = \text{abs}(\tanh(\Theta_{\text{time}}))$$

Same thing for x and calculate total reward :

$$\text{total\_reward} = 0.3 * x_{\text{reward}} + 0.7 * \Theta_{\text{reward}}$$

# Reward Structure

## *Correlation Matrix*

The values of 0.3 and 0.7 were retrieved from correlation matrix calculated from the original reward structure.

Absolute values of x-related terms and theta-related terms were combined to find the percentage of the total.

Example:

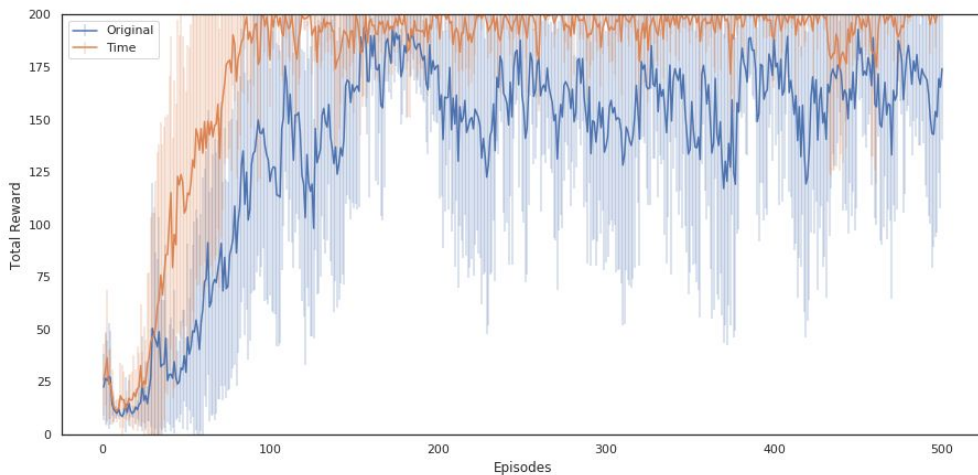
$$x\_coefficient = (0.56 + 0.05) / 2.06 \approx 0.3$$

	x	x_dot	theta	theta_dot	Total Reward
x	1.000000	0.685613	-0.209230	-0.524482	0.564305
x_dot	0.685613	1.000000	0.333420	0.004916	0.051151
theta	-0.209230	0.333420	1.000000	0.477523	-0.760483
theta_dot	-0.524482	0.004916	0.477523	1.000000	-0.682016
Total Reward	0.564305	0.051151	-0.760483	-0.682016	1.000000

## *Modified\_version2*

$$\text{Modified reward} = 0.3 * x\_r + 0.7 * \Theta\_r$$

Total Reward vs Episodes averaged over 10 runs

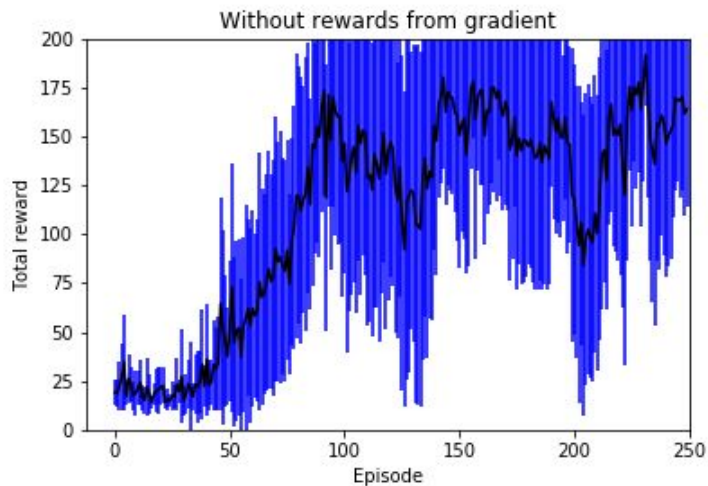




# Reward Structure

## *Original*

A reward of value 1.0 is granted so long as the pole angle is  $-41.8^\circ < \theta < 41.8^\circ$  and the cart position is  $-2.4 < x < 2.4$



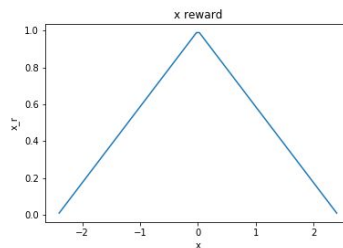
## *Modified*

The reward value is varied depending on how close the agent comes to failing during the step

*Calculating the Reward:*

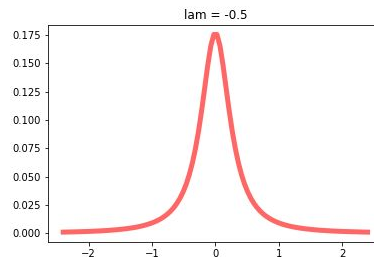
$\theta_r$  and  $x_r$ :

Triangular Function



$\theta_{dot_r}$  and  $x_{dot_r}$ :

Based on lam=-0.5 tukey  
lambda distribution

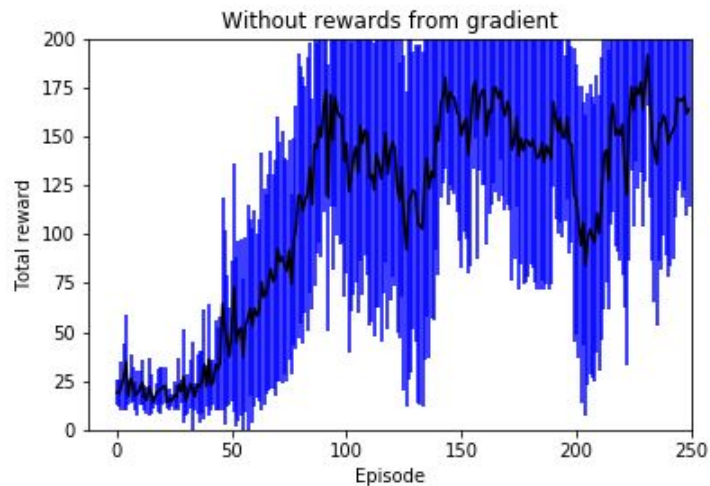


$$\text{Modified reward} = x * \theta + \theta_r * x_r$$

# Reward Structure

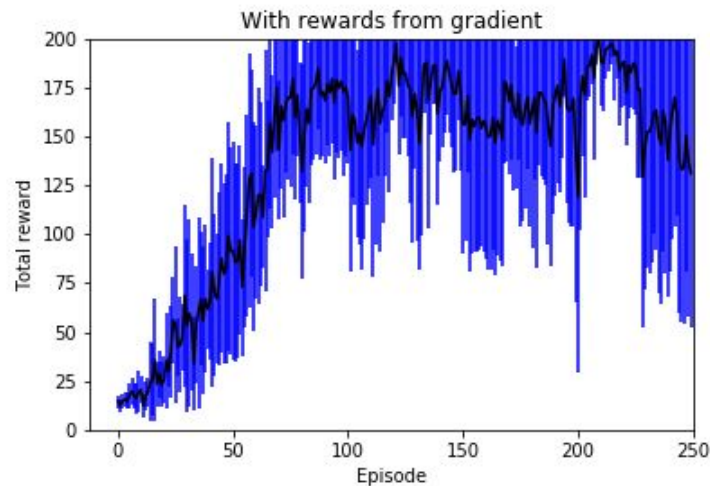
## *Original*

A reward of value 1.0 is granted so long as the pole angle is  $-41.8^\circ < \theta < 41.8^\circ$  and the cart position is  $-2.4 < x < 2.4$



## *Modified*

The reward value is varied depending on how close the agent comes to failing during the step



$$\text{Modified reward} = x * \theta * \theta_r * x_r$$

# Conclusion

- Convergence and stability of the model improved with modified reward structure
- Incorporation of physics into the neural network may improve the results significantly

