

# Simple RL Algorithm Comparisons

## 1. Q-Learning (Tabular)

Scenario: A robot exploring a 3-room house.

States: Kitchen, Living Room, Bathroom

Actions: Go Left, Go Right

Rewards: +1 in Kitchen, 0 elsewhere

You learn a table of expected rewards:  $Q[\text{state}][\text{action}] \rightarrow \text{expected reward}$

Example:  $Q[\text{Living Room}][\text{Go Left}] = 0.8$

Update Rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha * (r + \gamma * \max_{a'} Q(s', a') - Q(s, a))$$

Intuition: I thought this action was worth 0.2, but I got a reward of 1. I should increase my estimate!

## 2. DQN (Deep Q-Network)

Scenario: You're playing Atari (e.g., Breakout) with pixel inputs.

State: Image frame

Action: Discrete actions (left, right, no-op)

Model: A neural network predicting  $Q(s, a)$

You learn a function that maps state to Q-values for all actions:

$Q(s, a) \rightarrow \text{predicted reward}$

Update Rule:

$$\text{Loss} = (Q(s, a) - [r + \gamma * \max_{a'} Q_{\text{target}}(s', a')])^2$$

Intuition: My model said this action gives 0.4 reward, but it's really 0.8. Let's update the weights to fix that.

## 3. PPO (Proximal Policy Optimization)

Scenario: A robot balancing a pole, learning which direction to push.

State: Position, velocity, angle, etc.

# Simple RL Algorithm Comparisons

Action: Push left or right

Model: Policy network  $\pi(a | s)$  + Value network

You learn a probability distribution over actions:

$\pi(a | s)$  -> action probabilities

Example:  $\pi(\text{push-left} | \text{state}) = 0.6$

Update Rule:

$$L = E[\min(r_t * A_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) * A_t)]$$

Where  $r_t = \pi(a_t | s_t) / \pi_{\text{old}}(a_t | s_t)$

Intuition: This action was better than expected - increase its probability, but not too much!

## Summary Table

Q-Learning:

- Learns: Q-table (expected reward per action)
- Analogy: Robot exploring rooms
- Update: Adjust estimate if actual reward is different from expectation

DQN:

- Learns: Neural network for  $Q(s, a)$
- Analogy: Atari agent
- Update: Predict future reward using a neural net

PPO:

- Learns: A policy (probabilities over actions)
- Analogy: Pole-balancing robot
- Update: Slightly adjust the policy based on action advantage