

MC Control
SARSA ✓
Q-Learning }
Function Approximation^{TD} for V_{π}

Prof. Subrahmanya Swamy

Model-Free Setting so far..

• Prediction (To find V_π for a given π)

- MC First-Visit ✓
- MC Every-Visit ✓
- TD ✓
- N-step TD ✓
- MC Off-Policy ✓ π "b"

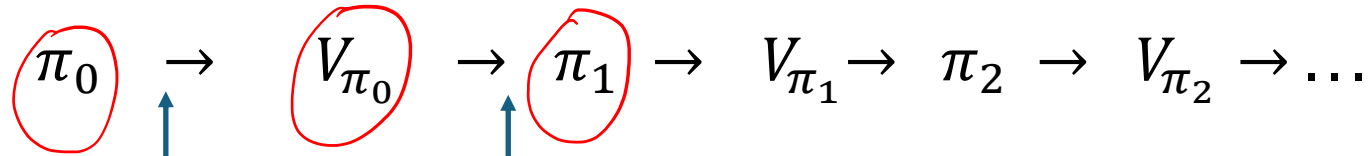
• "Control" (To find π^* of an MDP)

GPI. J.P.E.

- GPI with MC
- GPI with TD

Model Known: Policy Iteration (PI)

Policy Iteration



Repeat till

$$\pi_{k+1} = \pi_k$$



$$\pi_k = \pi_*$$

- Policy Evaluation: $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$

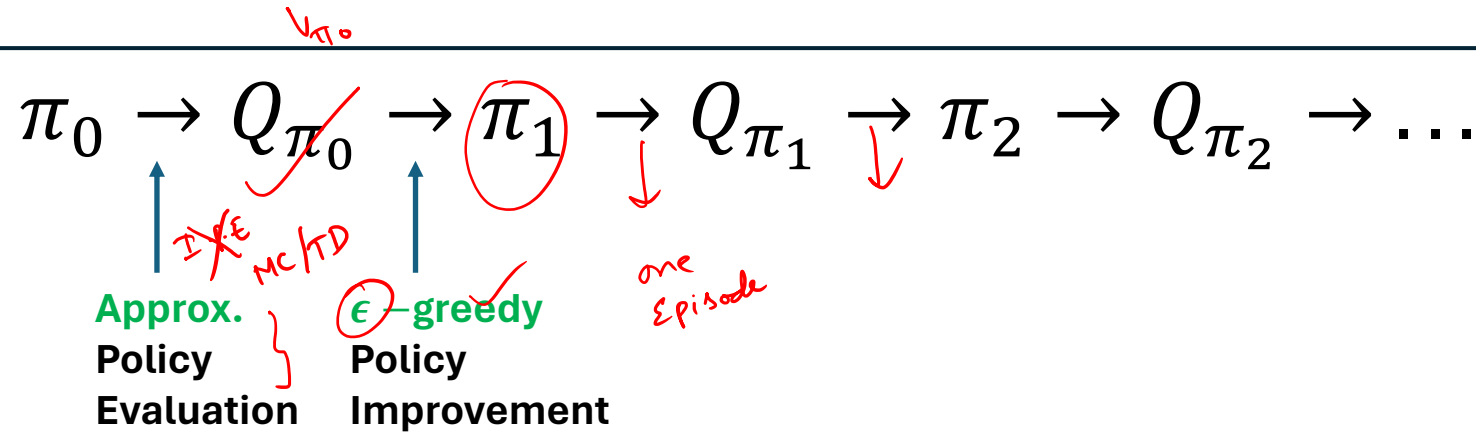
Repeat till $V_{k+1} = V_k$

- Policy Improvement: $\pi_{i+1}(s) := \operatorname{argmax}_a R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s')$

Policy Evaluation is stopped when $V_{k+1} \approx V_k$

Model Unknown: Generalized Policy Iteration (GPI)

PI \rightarrow GPI

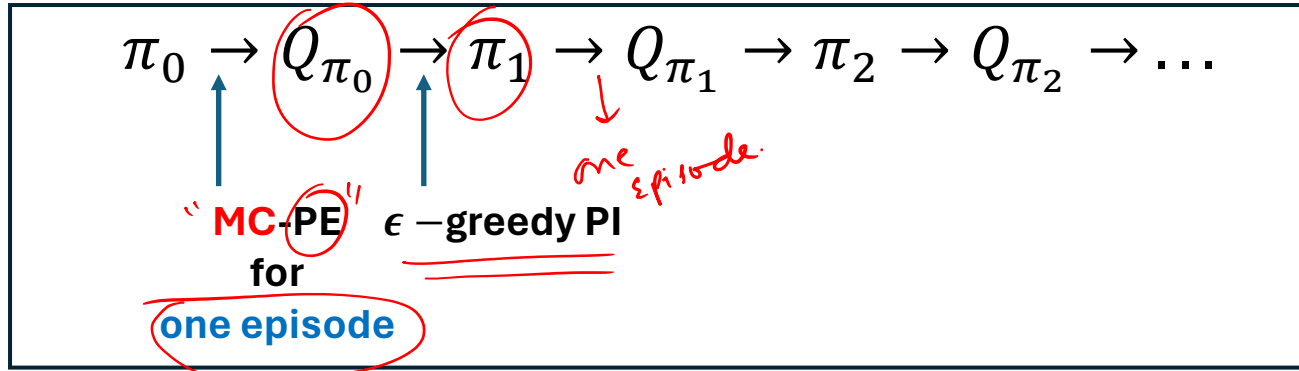


1. Q_{π_0} is evaluated
2. ϵ -greedy PI "Q(s,a)"
3. I.P.E is replaced with MC/TD.
 $V_{k+1} \approx V_k$.

When to stop the Approx. Policy Evaluation?

- MC: After one episode ✓
- TD: After one time-step ✓

MC (every visit) GPI : Pseudo Code

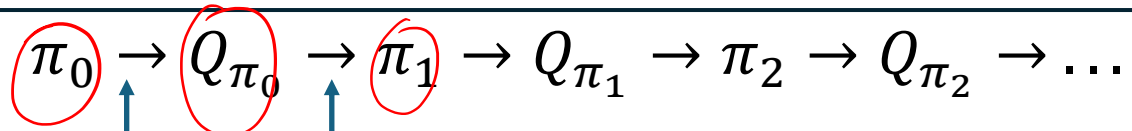


$$Q(s_0, A_0) \approx G_t$$

$$Q_{\text{new}}(s_0, A_0) = Q_{\text{old}}(s_0, A_0) + \alpha [G_t - Q_{\text{old}}(s_0, A_0)]$$

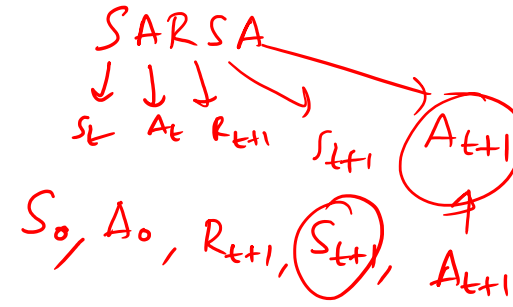
- Initialize $Q(s, a) = 0, \forall (s, a)$ ✓ Q_{π_0}
- Repeat** for each *episode*:
 - $\pi(s) = \epsilon$ -greedy w.r.t. $Q(s, a)$ ✓
 - Generate an episode following $\pi: (S_0, A_0), R_1, (S_1, A_1), R_2, (S_2, A_2), \dots, S_T$
 - Repeat** for each *time-step* t in the episode:
 - Compute $G_t = \sum_{i=t+1}^T \gamma^{i-t-1} R_i$
 - Update $Q(S_t, A_t) = Q(S_t, A_t) + \alpha (G_t - Q(S_t, A_t))$ ✓
- Output:** $\pi^* = \text{greedy}(Q)$

SARSA for π^* : Pseudo Code



TD-PE for "one time-step"
 ϵ -greedy PI

GPT



- Initialize $Q(s, a) = 0, \forall (s, a)$ ✓
- **Repeat** for each "episode":
 - Initialize S_0 randomly ✓
 - Sample $A_0 \sim \epsilon$ -greedy w.r.t. $Q(S_0, a)$ ✓
 - **Repeat** for each time-step t in the episode: ✓
 - Take action A_t and observe R_{t+1} and S_{t+1} ✓
 - Sample action $A_{t+1} \sim \epsilon$ -greedy w.r.t. $Q(S_{t+1}, a)$ ✓
 - $Q(S_t, A_t) = Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$ ✓
- **Output:** $\pi^* = \text{greedy}(Q)$ ✓

What ϵ and α to use?

$$\pi^*$$

ϵ

\rightarrow

" k^{th} Episode"

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \underline{Q^*(s, a)}$$

greedy w.r.t Q .

$$\textcircled{1} \quad \underline{\hat{\epsilon}_k} = \frac{1}{k} \rightarrow 0$$

$$\textcircled{2} \quad \sum_k \alpha_k = \infty, \quad \sum_k \alpha_k^2 < \infty. \checkmark$$

$$\alpha_k = \frac{1}{k}$$

$$\sum_k \frac{1}{k} = \infty$$

$$\sum_k \frac{1}{k^2} < \infty.$$

$$\epsilon_k = 1/k \checkmark$$

$$\alpha_k = 1/k \checkmark$$

MC
SARSA
TD - GPI
Q-learning

GPI.

Q-Learning

- **SARSA**: Based on Policy Iteration
- **Q-Learning**: Based on Value Iteration (Asynchronous)
- Value Iteration:

PI \rightarrow Q PI.

VI

$$\begin{aligned} \checkmark Q^*(s, a) &= R_s^a + \gamma \sum_{s'} P_{ss'}^a \underline{V^*(s')} & V^*(s) &= \max_a Q^*(s, a) \\ &= \checkmark R_s^a + \gamma \sum_{s'} \checkmark P_{ss'}^a (\max_{a'} Q^*(s', a')) \end{aligned}$$

$$\begin{aligned} \Rightarrow Q_{k+1}(s, a) &= R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q_k(s', a') \\ &\stackrel{Q_k(s, a)}{=} \text{repeat till.} \quad Q_{k+1} = Q_k. \\ &\quad \downarrow \\ &\quad \underset{(s, a) \checkmark}{\operatorname{argmax}} \left(Q_k(s, a) - Q_{k+1}(s, a) \right) \\ &\quad \underset{s}{\operatorname{argmax}} \left(V_{k+1}(s) - V_k(s) \right) \end{aligned}$$

Q-Learning

Value Iteration:

- **SARSA**: Based on Policy Iteration
- **Q-Learning**: Based on Value Iteration (Asynchronous)
- **Value Iteration**:

$$\begin{aligned} \checkmark Q^*(s, a) &= R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s') \\ &= R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q^*(s', a') \end{aligned}$$

$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

- **Q-Learning**:

$$Q_{new}(S_t, A_t) = Q_{old}(S_t, A_t) + \alpha (\checkmark R_{t+1} + \gamma \max_{a'} Q_{old}(S_{t+1}, a') - Q_{old}(S_t, A_t))$$

Which state and action pair should be updated in Q-learning? (s, a)

- Random policy ✓
 - ϵ -greedy w.r.t Q ✓
- $(S_t, A_t) \rightarrow \dots$

Q-Learning: Pseudo Code

- Initialize $Q(s, a) = 0, \forall(s, a)$ ✓
- **Repeat** for each episode:
 - Initialize S_0 randomly ✓
 - **Repeat** for each time-step t in the episode:
 - Sample action $A_t \sim \epsilon$ -greedy w.r.t. $Q(S_t, a)$ ✓
 - Take action A_t and observe R_{t+1} and S_{t+1} ✓
 - ✓ $Q(S_t, A_t) = Q(\underline{S_t}, \underline{A_t}) + \alpha (\underline{R_{t+1}} + \gamma (\max_a Q(\underline{S_{t+1}}, \underline{a})) - Q(S_t, A_t))$ ✓
- **Output:** π^* = "greedy (Q)" ✓

off-policy

Q^*

SARSA Vs Q-Learning

• SARSA

- On-Policy ✓
- Based on Policy Iteration ✓
- Converges to the best among ϵ -soft policies if fixed ϵ is chosen ✓
- Converges to π^* if ϵ is decreased to zero with time ✓

" ϵ -greedy"

Optimal π^* .

" $\epsilon = 0$ "

$$= \epsilon_k = \frac{1}{k}$$

Q^k

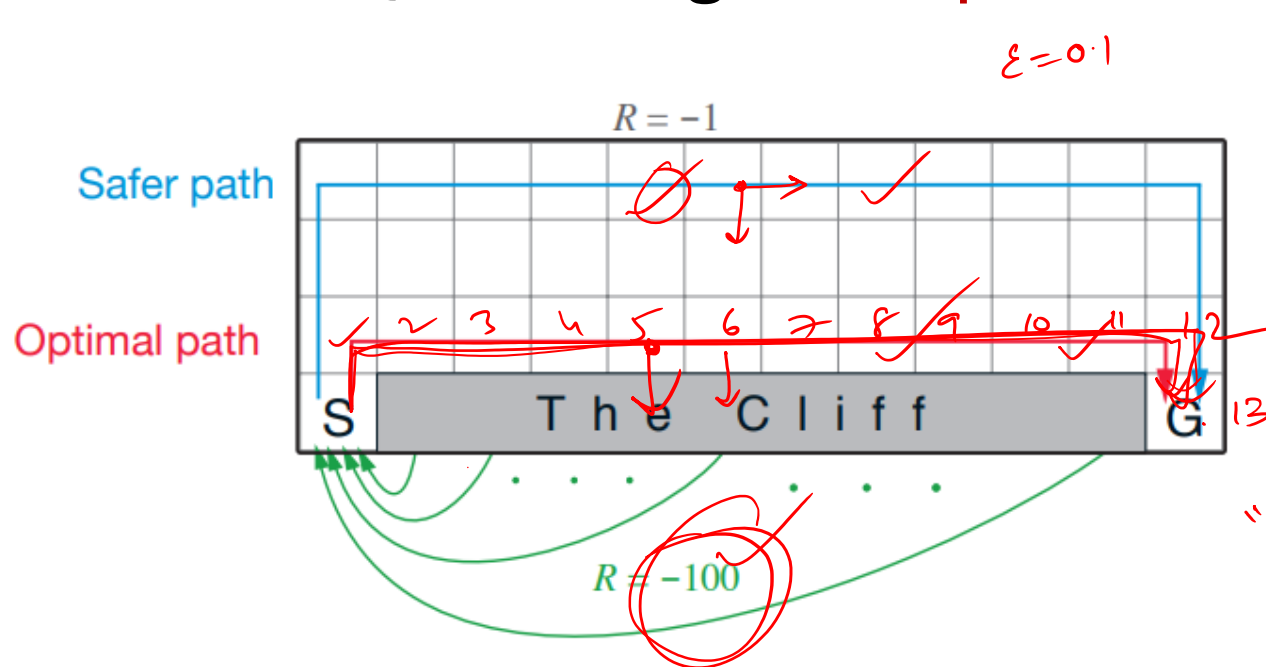
ϵ -greedy

(s_t, A_t)

• Q-Learning

- Off-Policy ✓
- Based on Value Iteration ✓
- Converges to π^* even for fixed $\epsilon = 0.1$

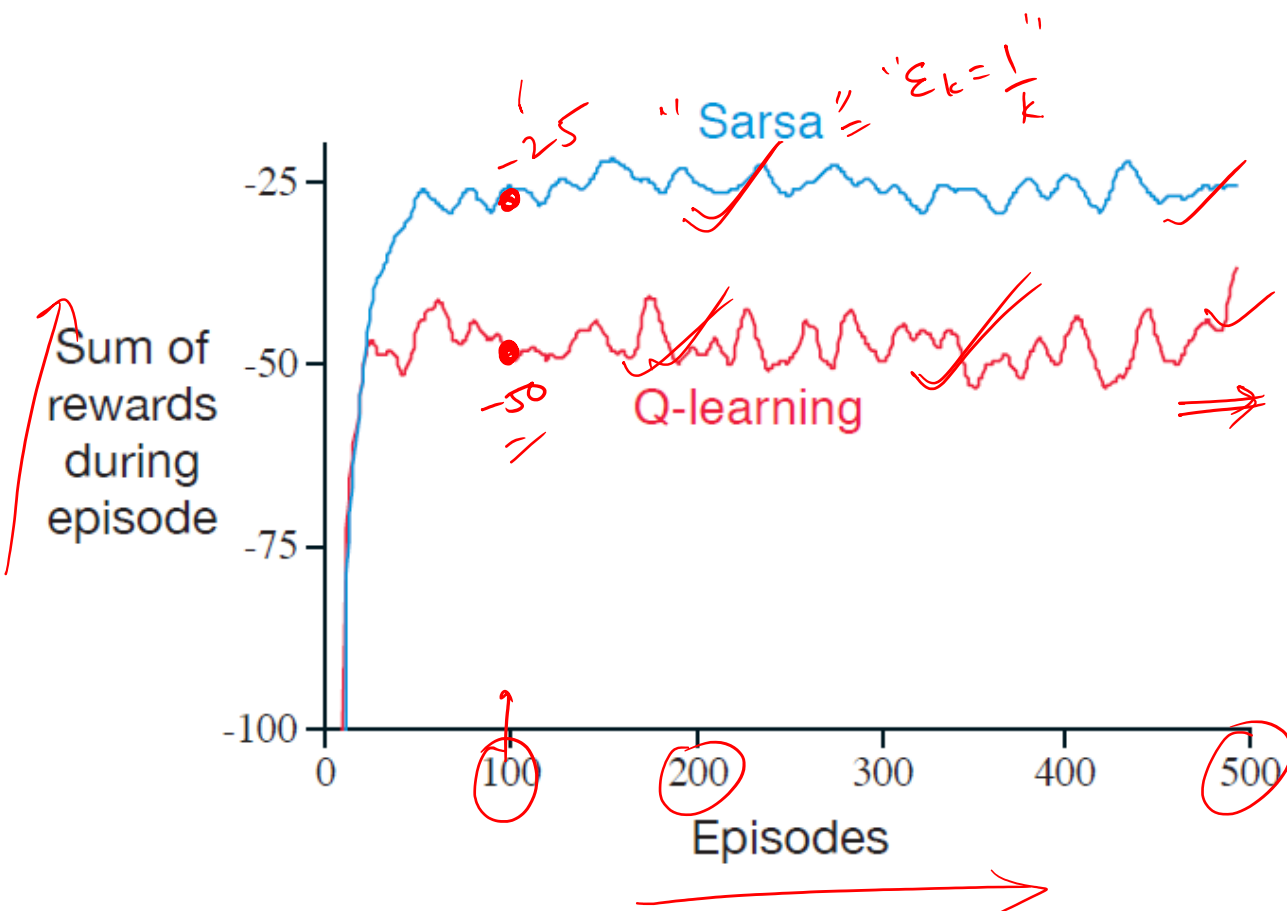
SARSA Vs Q-Learning: Example



- SARSA: Learns Safer path
- Q-Learning: Learns Optimal path

- **Aim**: To go in the shortest path from the Start state to the Goal state
- Reward of -100 for transition into the **Cliff** region
- Reward of -1 for every **other transition**

SARSA Vs Q-Learning: Example



- ~~SARSA~~: Learns Safer path
- ~~Q-Learning~~: Learns Optimal path

- Why does Q-learning have a worse reward, although it learned the optimal policy?

Function Approximation for Large State spaces

Prof. Subrahmanya Swamy

Large State Spaces in Real-time Applications

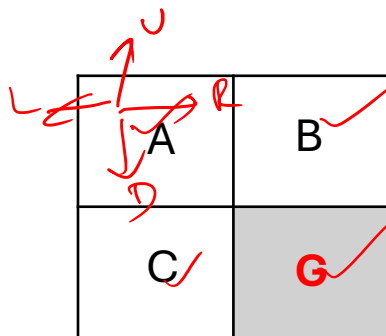
Go game



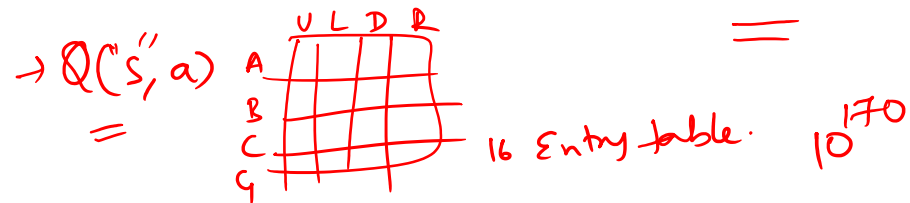
$\sim 10^{170}$ states $>$ # Atoms in the universe

Tabular Methods

- **Small** state/action space: $Q(s, a)$ Table maintained for each (s, a) explicitly



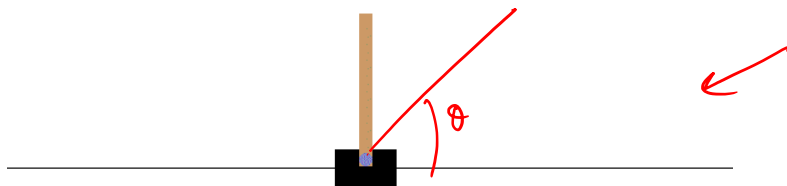
Q-Table has 4 states x 4 actions = 16 entries



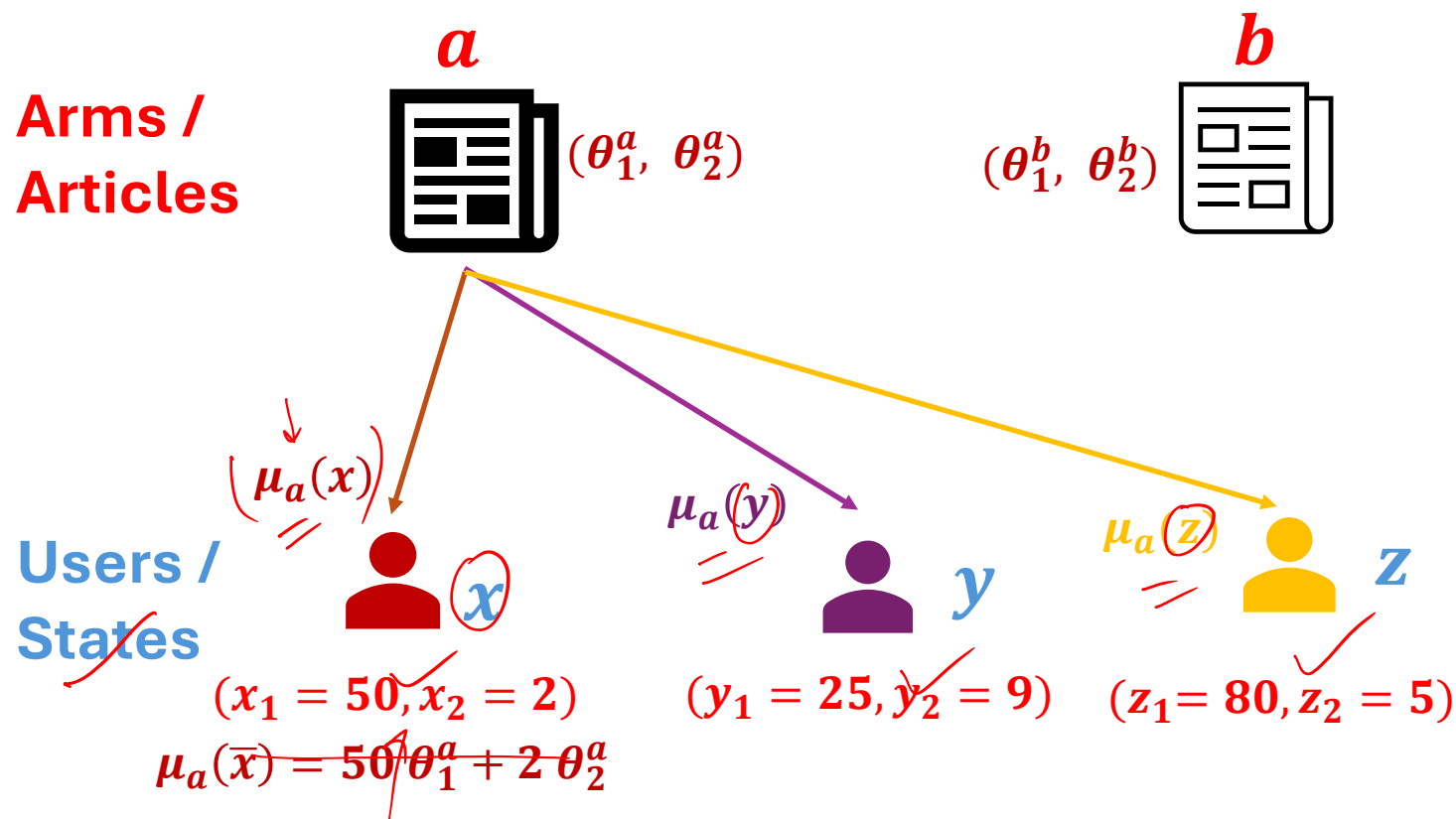
- **Large** state/action spaces: **Not feasible!**

- **Continuous** state/action spaces: **Not feasible!**

In Approx
Contextual Bandits



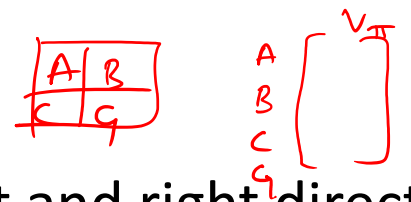
Bandits + Supervised Learning



- User (state) represented by features such as age, income
 $\bar{x} = (x_1, x_2)$
- Model the expected reward for user \bar{x} for pulling arm **a** as
 $\mu_a(\bar{x}) = \theta_1^a x_1 + \theta_2^a x_2$

Features and Function Approximation

features



Cartpole: The goal is to balance the pole by applying forces in the left and right direction

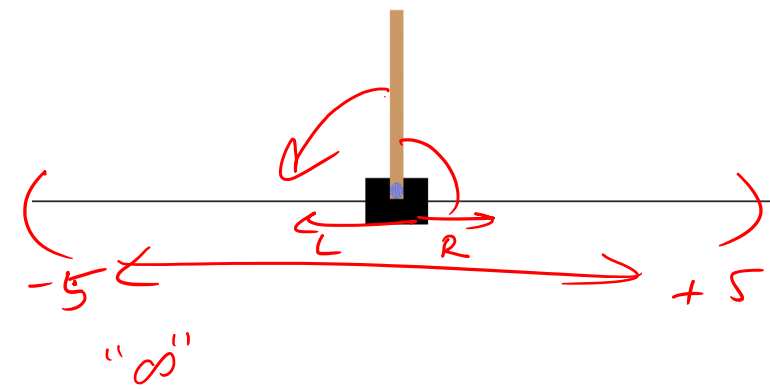
State Features

$s = (\check{s}_1, \check{s}_2, \check{s}_3, \check{s}_4)$

$V_{\pi}(s) \leftarrow$

$Q(s, a)$

| State \bar{s} | Min | Max |
|-----------------------------|--------------------------|------------------------|
| Cart Position s_1 | -4.8 | 4.8 |
| Cart Velocity s_2 | -Inf | Inf |
| Pole Angle s_3 | ~ -0.418 rad (-24°) | ~ 0.418 rad (24°) |
| Pole Angular Velocity s_4 | -Inf | Inf |



Value fn Approx

$V_{\theta}(s) \approx s_1\theta_1 + s_2\theta_2 + s_3\theta_3 + s_4\theta_4$
 $= Q(\bar{s}, a)$

a Action Features

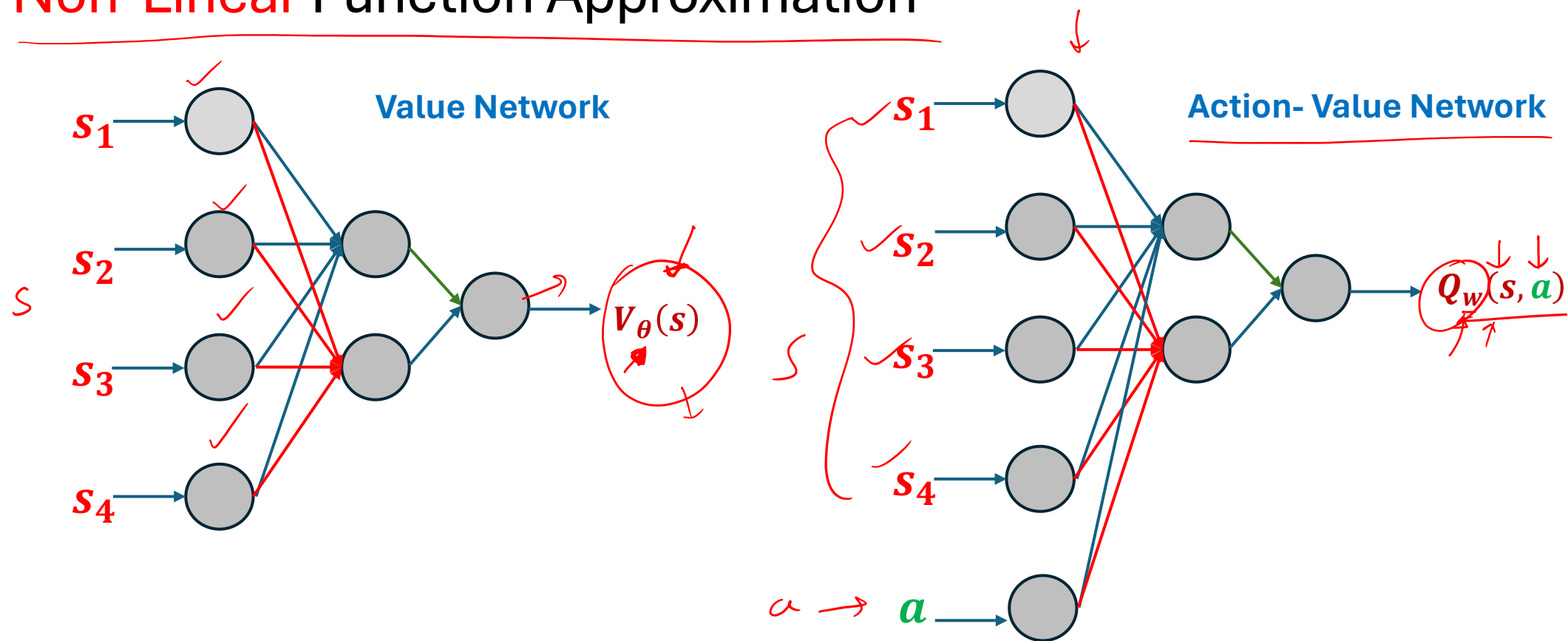
- 0: Push the cart to the LEFT
- 1: Push the cart to the RIGHT

Linear "w"

Q fn Approx

$Q_w(s, a) \approx s_1w_1 + s_2w_2 + s_3w_3 + s_4w_4 + aw_5$

Non-Linear Function Approximation



Neural Network-based Function Approximation

How to find the weights of the neural network?

Function approx. for V_π π

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_\pi \left[\left(\underline{V_\pi(s)} - \underline{V_\theta(s)} \right)^2 \right]$$

$$\begin{pmatrix} \hat{s}_1 \rightarrow \check{V_\pi(s_1)} \\ \check{s}_2 \rightarrow \check{V_\pi(s_2)} \end{pmatrix} \leftarrow s_3$$

$$\theta_{new} = \theta_{old} + \alpha \cdot \frac{\partial L}{\partial \theta}$$

Stochastic Gradient descent:

$$\theta_{new} = \theta_{old} + 2 \alpha (\underline{V_\pi(s)} - V_\theta(s)) \nabla V_\theta(s)$$

Challenge:

RL

- " $V_\pi(s)$ unknown"
- No training data available ✓

MC Function approx. for V_π

MC/TD

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_\pi \left[\left(V_\pi(s) - V_\theta(s) \right)^2 \right]$$

$$\theta_{\text{new}} = \theta_{\text{old}} + 2 \alpha \left(V_\pi(s) - V_\theta(s) \right) \nabla V_\theta(s)$$

SGD

$V_\pi(s) \approx \hat{G}_t$ starting from state s

$$\theta_{\text{new}} = \theta_{\text{old}} + 2 \alpha \left(\hat{G}_t - V_\theta(s) \right) \nabla V_\theta(s) \quad \checkmark$$

$$V_\pi(s) = \mathbb{E}_\pi \left[G_t \mid S_t = s \right]$$

$s \xrightarrow{\quad} s_\tau \rightarrow \hat{G}_t$

TD Function approx. for V_π

"TD"

π Tab

$$\theta_{\text{new}} = \theta_{\text{old}} + 2\alpha \left(\underbrace{V_\pi(s)}_{\text{MC}} - V_\theta(s) \right) \nabla V_\theta \quad \text{SGD}$$

\downarrow
 $R_{t+1} + \gamma V_\theta(s_{t+1})$

$$\theta_\pi = \theta^*$$

MC: $\underline{G_t}$

$$\begin{aligned} V_\pi(s) &= \mathbb{E}_\pi [G_t \mid s_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} \mid s_t = s] \end{aligned}$$

$$\approx R_{t+1} + \gamma \underline{V_\pi(s_{t+1})}$$

$$\approx \boxed{R_{t+1} + \gamma V_\theta(s_{t+1})}$$