# Introduction To Reinforcement Learning

Thien-Minh Nguyen, PhD

Centre for Advance Robotics Technology Innovation

Feb 2025

1

# Outline

- Recent RL applications in robotics.
- Overview of RL
- Basic Methods to solve the RL problems
- Tutorials
  - Tic-tac-toe. Complete MDP with *Value Iteration Technique*.
  - Cartpole. Small scaled complete RL problem → benchmarking and analysis

# Objectives

- Exposure to mathematical formulism of RL.
- Familiarize with <u>basic concepts</u> of Reinforcement Learning (RL).

  In the context of RL...

  - Agent, environment, observations, state, reward, action, value, return, discount ...
  - Evaluation, Iteration, Improvement, Value Iteration ...
  - Monte Carlo, Off-policy
  - Temporal Difference, Q-learning, Sarsa
  - Function Approximation
  - Policy Gradient Methods

# Why RL ain't work?

- Sample Inefficient
- Can be solved by other methods
- Always requires a reward function
- Reward function design is difficult
- Local optima hard to escape
- Overfitting
- Unstable and hard to reproduce

# Why RL *works* now?

- Sample Inefficient → Cost of experiment↓



$75,000

Pricing considerations

The Spot robot dog price isn't budget-friendly: Spot starts at $75,000. That's before adding payloads, sensors, or software packages. 3 Mar 2025

Standard Bots
https://standardbots.com › blog › spot-robot

The Spot Robot by Boston Dynamics: Features & Use Cases



Welcome to Isaac Lab!

Isaac Lab Documentation

**Version**
main

IsaacSim 4.5.0    Stars 3.1k

Search... ctrl K

**Getting Started**

Isaac Lab Ecosystem
Local Installation
Running Isaac Lab in the Cloud

**Overview**

Developer's Guide
Core Concepts
Available Environments
Reinforcement Lear

**Isaac Lab** is a unified and modular framework for robot learning that aims to simplify common workflows in

The ANYmal robot, a four-legged robot made by ANYbotics for industrial customers, costs around $150,000 and includes the full autonomy platform with LIDAR and a docking station, but excludes payloads, self-charging docks, and autonomous capabilities.

## GENESIS

latest    User Guide    API Reference    Roadmap

**Section Navigation**

**Overview**
💡 What is Genesis
💎 Why A New Physics Simulator
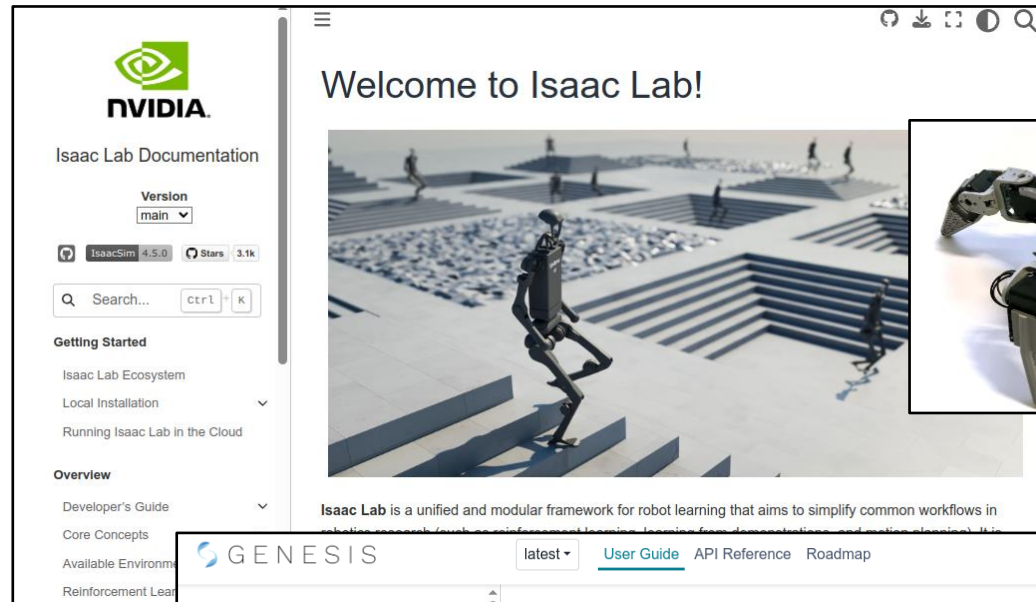🔧 Installation
🎯 Genesis Vision & Mission

**Getting Started**
👋 Hello, Genesis
🎨 Visualization & Rendering
🎮 Control Your Robot
🚀 Parallel Simulation
🦿 Inverse Kinematics & Motion Planning
🤖 Advanced and Parallel IK
🧊 Beyond Rigid Bodies
📡 Interactive Information Access and

User Guide › Training Locomotion Policies with RL

## 🦾 Training Locomotion Policies with RL

Genesis supports parallel simulation, making it ideal for training reinforcement learning (RL) locomotion policies efficiently. In this tutorial, we will walk you through a complete training e obtaining a basic locomotion policy that enables a Unitree Go2 Robot to walk.

This is a simple and minimal example that demonstrates a very basic RL training pipeline in and with the following example you will be able to obtain a quadruped locomotion policy tha deployable to a real robot very quickly.

**Note**: This is *NOT* a comprehensive locomotion policy training pipeline. It uses simplified re terms to get you started easily, and does not exploit Genesis's speed on big batchsizes, so serves basic demonstration purposes.

**Acknowledgement**: This tutorial is inspired by and builds several core concepts from Legged Gym.

## UNITREE

Home    Professional    Education & Industry    Humanoid Robot    Accessories    Contact    Disclaimer

Home › Unitree Go1

### Unitree Go1

$2,700.00 USD

📦 Ready to ship

🚚 Shipping costs $1000 per unit

$ Duty not included, US needs to charge 25 percent duty

💬 Contact sales for Go1 Edu price

**Version**

Go1 Air    Go1 pro

Chat with us 👋

6

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment** ↓

- Some problems can be solved by other methods
  → and many others can be solved by RL



NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

⌂ > Nanyang Business School > News & Events > News

Published on 14 Nov 2023

**NBS Knowledge Lab Interdisciplinary Distinguished Speaker Series webinar: Reinforcement Learning for Quantitative Trading**



📖 README  ⚖ MIT license

🔗 **DeepSeek-R1**

deepseek



nature

Explore content ⌄     About the journal ⌄     Publish with us ⌄

nature > articles > article

Article | Open access | Published: 05 October 2022

**Discovering faster matrix multiplication algorithms with reinforcement learning**

AlphaGo

AlphaGo mastered the ancient game of Go, defeated a Go world champion, and inspired a new era of AI systems.

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment** ↓

- Some problems can be solved by other methods
  → and many others can be solved by RL



Learning Dynamic Weight Adjustment for Spatial-Temporal Trajectory Planning in Crowd Navigation

Muqing Cao*, Xinhang Xu*, Yizhuo Yang*, Jianping Li, Tongxing Jin, Pengfei Wang, Tzu-Yi Hung, Guosheng Lin, and Lihua Xie[1] *Fellow, IEEE*



Reinforcement Learning Meets Visual Odometry

Nico Messikommer*, Giovanni Cioffi*,
Mathias Gehrig, and Davide Scaramuzza

Dept. of Informatics, University of Zurich
{nmessi,cioffi,mgehrig,sdavide}@ifi.uzh.ch

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment** ↓

- Some problems can be solved by other methods
  → and many others can be solved by RL

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment ↓**
- Some problems can be solved by other methods → **Some can be solved by RL**
- Always requires a reward function
- Reward function design is difficult
- Local optima hard to escape
- Overfitting
- Unstable and hard to reproduce
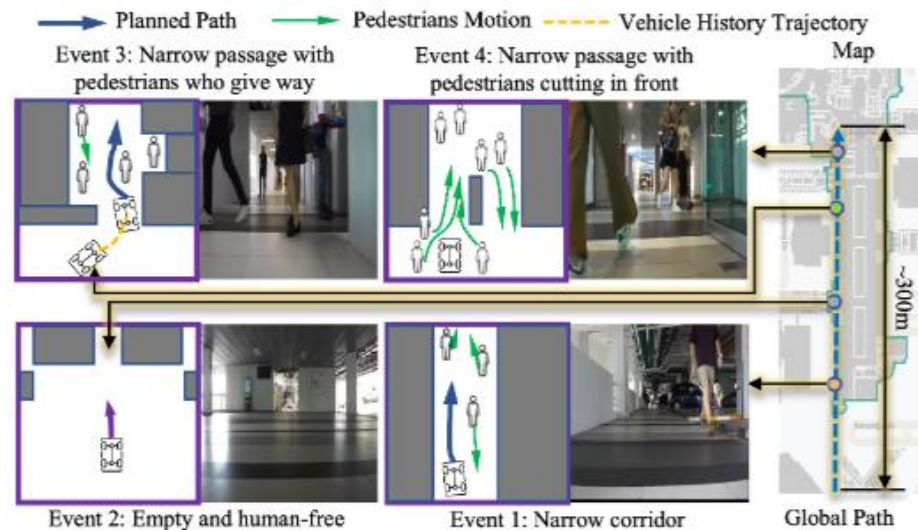
Fermat's equation:

$$x^n + y^n = z^n$$

This equation has no solutions in integers for $n \geq 3$.

## Curriculum-Based Reinforcement Learning for Quadrupedal Jumping: A Reference-free Design

Vassil Atanassov*, Jiatao Ding*, Jens Kober, Ioannis Havoutis, Cosimo Della Santina

REWARDS DEFINITION. THE LIGHT ORANGE COLOUR INDICATES TASK-BASED REWARDS, WHILE THE LIGHT PURPLE SHADE DESCRIBES REGULARISATION REWARDS. $w_\times$ IS THE WEIGHT, $\sigma_\times$ IS A SCALING FACTOR FOR THE EXPONENTIAL KERNEL, $e(\cdot)$ AND $\log(\cdot)$ SEPARATELY DENOTE THE EXPONENT AND LOGARITHM OPERATION.

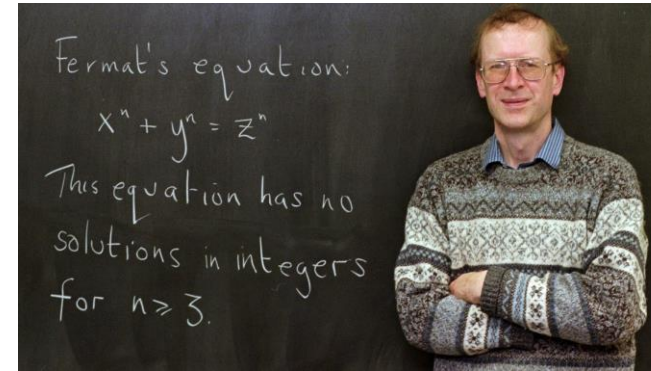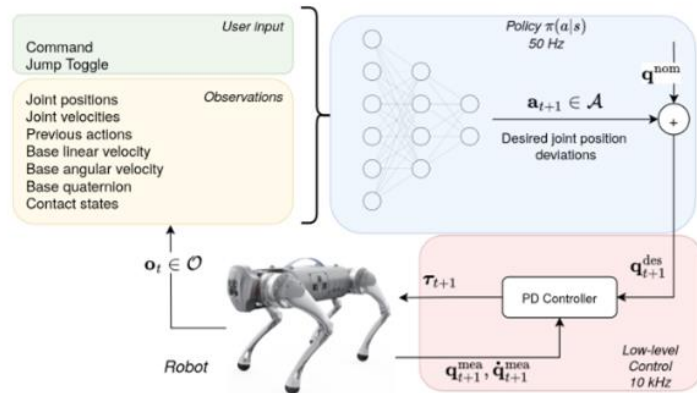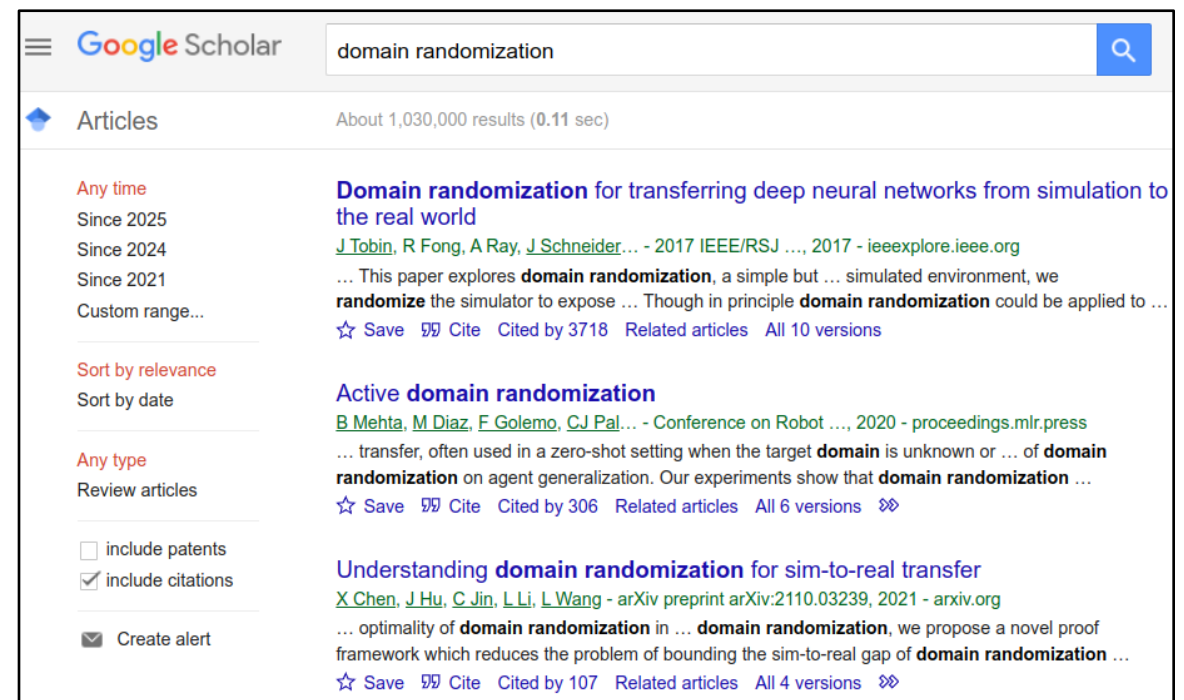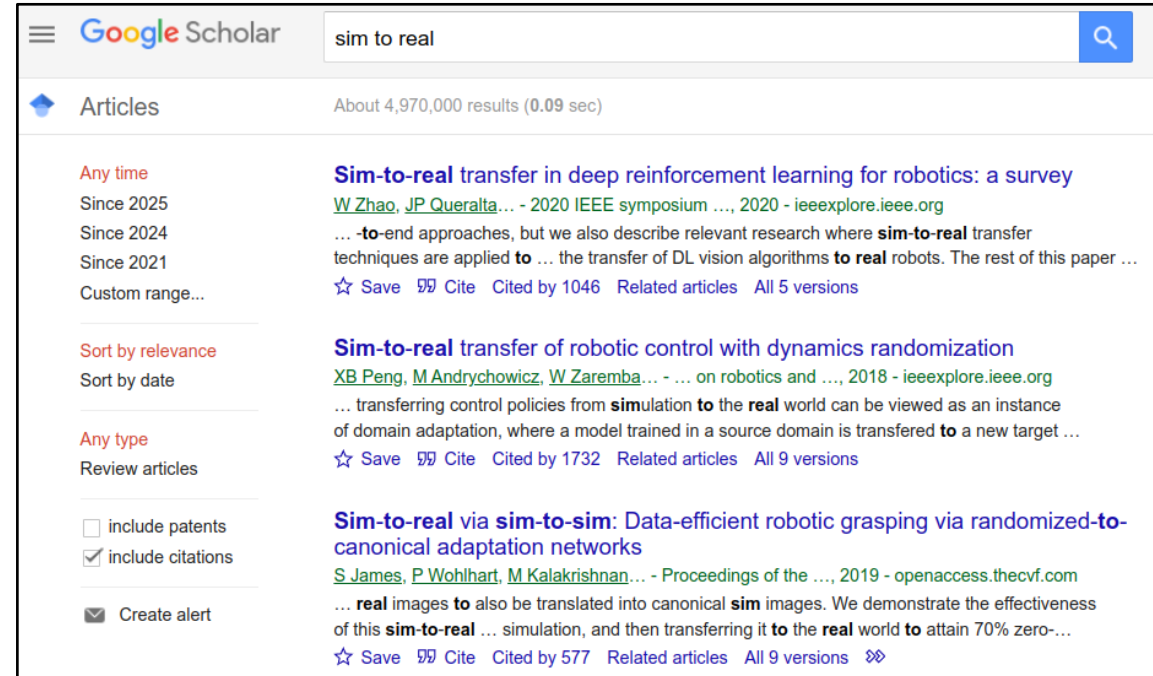| Name | Type | Stance | Flight | Landing |
|---|---|---|---|---|
| Landing position | Single | 0 | 0 | $w_{\mathbf{p}}(e(-\sum \|\mathbf{P}_{\mathrm{land}} - \mathbf{p}_{\mathrm{des}}\|^2)/\sigma_{p,\mathrm{land}})$ |
| Landing orientation | Single | 0 | 0 | $w_{\mathrm{ori}}(e(-\|\log(\bar{\mathbf{q}}_{\mathrm{land}}^{-1} * \bar{\mathbf{q}}_{\mathrm{des}})\|^2)/\sigma_{\mathrm{ori,land}})$ |
| Max height | Single | 0 | 0 | $w_h(e(\|h_{\max} - 0.9\|^2)/\sigma_{p_z,\max})$ |
| Jumping | Single | 0 | 0 | $w_{\mathrm{jump}}$ |
| Base Position | Continuous | $w_{p_z,\mathrm{st}}(e(-\|p_z - 0.20\|^2/\sigma_{p_z,\mathrm{st}}))$ | $w_{p_z,\mathrm{fl}}(e(-\|p_z - 0.7\|^2/\sigma_{p_z,\mathrm{fl}}))$ | $w_{\mathbf{p},\mathrm{l}}(e(-\sum \|\mathbf{p} - \mathbf{p}_{\mathrm{des}}\|^2/\sigma_{p,\mathrm{l}}))$ |
| Orientation Tracking | Continuous | $w_{\mathrm{ori,st}}(e(-\|\log(\bar{\mathbf{q}}_{\mathrm{base}}^{-1} * \bar{\mathbf{q}}_{\mathrm{des}})\|^2/\sigma_{\mathrm{ori,st}}))$ | 0 | $w_{\mathrm{ori,l}}(e(-\|\log(\bar{\mathbf{q}}_{\mathrm{base}}^{-1} * \bar{\mathbf{q}}_{\mathrm{des}})\|^2/\sigma_{\mathrm{ori,l}}))$ |
| Base linear velocity | Continuous | 0 | $w_{\mathbf{v}_{x,y}}(-e(\sum \|\mathbf{v}_{x,y} - \mathbf{v}_{\mathrm{des}}\|^2/\sigma_v))$ | 0 |
| Base angular velocity | Continuous | 0 | $w_{\boldsymbol{\omega}}(e(-\sum \|\boldsymbol{\omega} - \boldsymbol{\omega}_{\mathrm{des}}\|^2/\sigma_{\omega}))$ | $0.1 w_{\boldsymbol{\omega}}(e(-\sum \|\boldsymbol{\omega}\|^2/\sigma_{\omega}))$ |
| Feet clearance | Continuous | 0 | $w_{\mathrm{feet}}(\|p_{\mathrm{feet}} - p_{\mathrm{feet}}^0 + [0.0, 0.0, -0.15]\|^2)$ | 0 |
| Symmetry | Continuous | $w_{\mathrm{sym}}(\sum_{\mathrm{joint}} \|\mathbf{q}_{\mathrm{left}} - \mathbf{q}_{\mathrm{right}}\|^2)$ | | |
| Nominal pose | Continuous | $w_{\mathbf{q}}(e(-\sum_{\mathrm{joint}} \|\mathbf{q}_j - \mathbf{q}_{j,\mathrm{nom}}\|^2/\sigma_q)$ | $0.1 w_{\mathbf{q}}(e(-\sum_{\mathrm{joint}} \|\mathbf{q}_j - \mathbf{q}_{j,\mathrm{nom}}\|^2/\sigma_q)$ | $w_{\mathbf{q}}(e(-\sum_{\mathrm{joint}} \|\mathbf{q}_j - \mathbf{q}_{j,\mathrm{nom}}\|^2/\sigma_q)$ |
| Energy | Continuous | $w_{\mathrm{energy}}(\boldsymbol{\tau}^T \dot{\mathbf{q}})$ | | |
| Base acceleration | Continuous | $w_{\mathrm{acc}}\|\dot{\mathbf{v}}\|^2$ | | |
| Contact change | Continuous | $w_c \sum_{\mathrm{feet}} (c_{\mathrm{foot}}(t) - c_{\mathrm{foot}}(t-1))$ | | |
| Maintain Contact | Continuous | $w_{\mathrm{contact}} \sum_{\mathrm{feet}} c_{\mathrm{foot}}(t)$ | 0 | 0 |
| Contact forces | Continuous | $w_{F_c} \sum_{i=0}^{n_f} \|F_i - \bar{F}\|$ | | |
| Action rate | Continuous | $w_a \sum_{\mathrm{joint}} \|\mathbf{a}(t) - \mathbf{a}(t-1)\|^2$ | | |
| Joint acceleration | Continuous | $w_{\ddot{\mathbf{q}}} \sum_{\mathrm{joint}} \|\ddot{\mathbf{q}}_j\|^2$ | | |
| Joint limits | Continuous | $w_{q_{lim}} \sum_{\mathrm{joint}} \|\mathbf{q}_j - \mathbf{q}_{j,lim}\|^2$ | | |

Command
Jump Toggle

User input

Policy $\pi(a|s)$
50 Hz

$\mathbf{q}^{\mathrm{nom}}$

Joint positions
Joint velocities
Previous actions
Base linear velocity
Base angular velocity
Base quaternion
Contact states

Observations

$\mathbf{a}_{t+1} \in \mathcal{A}$

Desired joint position deviations

$\mathbf{o}_t \in \mathcal{O}$

$\boldsymbol{\tau}_{t+1}$

PD Controller

$\mathbf{q}_{t+1}^{\mathrm{des}}$

Low-level
Control
10 kHz

Robot

$\mathbf{q}_{t+1}^{\mathrm{mea}}, \dot{\mathbf{q}}_{t+1}^{\mathrm{mea}}$

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment ↓**

- Some problems can be solved by other methods → **Some can be solved by RL**

- Always requires a reward function

- Reward function design is difficult

- <span style="color:red">Local optima hard to escape</span>

- <span style="color:red">Overfitting</span>

- <span style="color:red">Unstable and hard to reproduce</span>

  → **Active research areas**

# Why RL *works* now?

- Sample Inefficient → **Cost of experiment** ↓

- Some problems can be solved by other methods → **Some can be solved by RL**

- Always requires a reward function

- Reward function design is difficult

- Local optima hard to escape

- Overfitting

- Unstable and hard to reproduce

→ **Active research areas**

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount ...

- Evaluation, Iteration, Improvement, Value Iteration ...

- Monte Carlo, Off-policy

- Temporal Difference, Q-learning, Sarsa

- Function Approximation

- Policy Gradient Methods

# Agent and Environment

- **Agent**: receives **observations** and **rewards**, generates **action**.

- **Environment**: receives **action**, produces **observation** and **reward**.



The robot belongs to which category?

*"All goals can be described by the maximization of expected cumulative reward"*

Code   Blame   148 lines (130 loc) · 4.48 KB
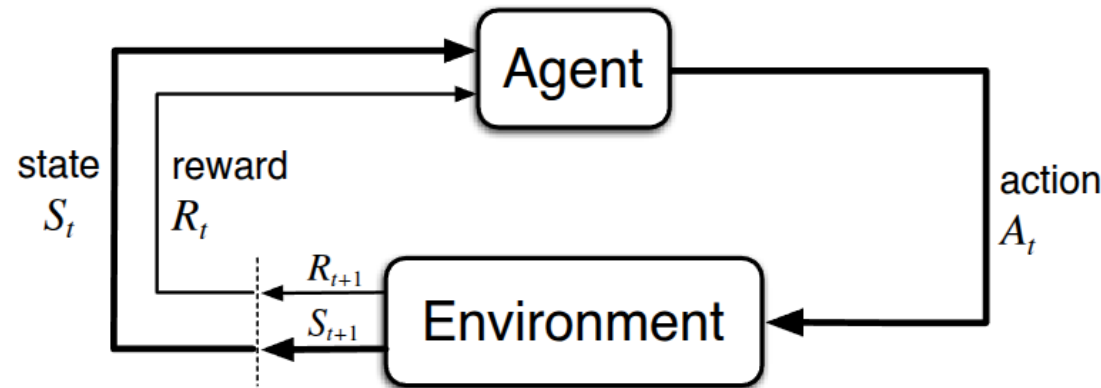
```python
53      class AnymalCFlatEnvCfg(DirectRLEnvCfg):
56          decimation = 4
57          action_scale = 0.5
58          action_space = 12
59          observation_space = 48
60          state_space = 0
61
62          # simulation
63  >       sim: SimulationCfg = SimulationCfg( ...
73          )
74  >       terrain = TerrainImporterCfg( ...
86          )
87
88          # scene
89          scene: InteractiveSceneCfg = InteractiveSceneCfg(num_envs=4096, env_spacing=4.0, replicate_physics=True)
90
91          # events
92          events: EventCfg = EventCfg()
93
94          # robot
95          robot: ArticulationCfg = ANYMAL_C_CFG.replace(prim_path="/World/envs/env_.*/Robot")
96          contact_sensor: ContactSensorCfg = ContactSensorCfg(
97              prim_path="/World/envs/env_.*/Robot/.*", history_length=3, update_period=0.005, track_air_time=True
98          )
99
100         # reward scales
101         lin_vel_reward_scale = 1.0
102         yaw_rate_reward_scale = 0.5
103         z_vel_reward_scale = -2.0
104         ang_vel_reward_scale = -0.05
105         joint_torque_reward_scale = -2.5e-5
106         joint_accel_reward_scale = -2.5e-7
107         action_rate_reward_scale = -0.01
108         feet_air_time_reward_scale = 0.5
109         undesired_contact_reward_scale = -1.0
110         flat_orientation_reward_scale = -5.0
```
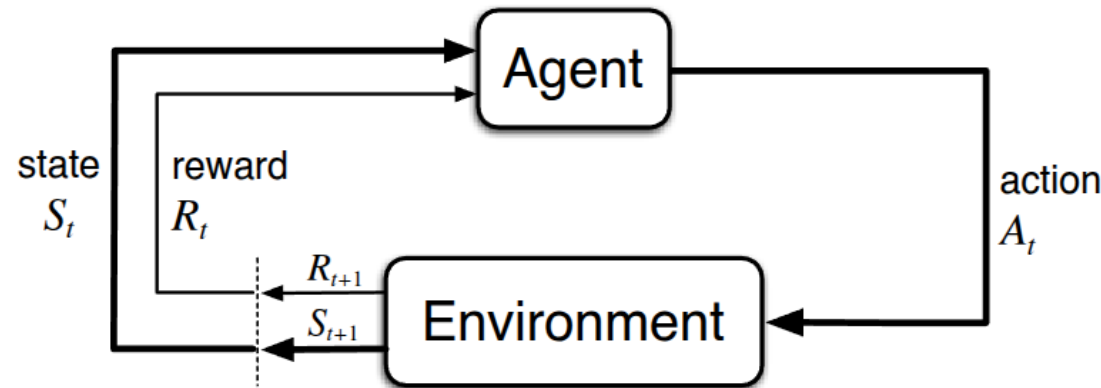
15

# Concepts (1)



## Definition

- A finite Markov Decision Process:
    - $R_t$ is the **reward**, **a scalar signal**
    - $A_t$ is the **action**, e.g., torque command, velocity command, chess moves ...
    - $S_t$ is the **state.**
    - $t \in \{0, 1, 2 ...\}, s \in \mathcal{S}, a \in \mathcal{A}(s), r \in \mathcal{R} \subset \mathbb{R}$
    - The dynamics between agent and environment is specified as:
    $$\mathcal{P} \triangleq p(s', r | s, a) = \text{Prob}(S_{t+1} = s', R_{t+1} | S_t = s, A_t = a)$$
- $H_t \triangleq (S_1, R_1, A_1 ... S_t, R_t, A_t)$, the **trajectory.**
- $O_t = h(H_t)$ is the **observation**, e.g., image, can be IMU reading, lidar scan ...

# Concepts (1)



## Definition

- A finite Markov Decision Process:
    - $R_t$ is the **reward**, **a scalar signal**
    - $A_t$ is the **action**, e.g., torque command, velocity command, chess moves …
    - $S_t$ is the **state.**
    - $t \in \{0, 1, 2 \ldots\}, s \in \mathcal{S}, a \in \mathcal{A}(s), r \in \mathcal{R} \subset \mathbb{R}$
    - The dynamics between agent and environment is summarized in:
$$\mathcal{P} \triangleq p(s', r | s, a) = \text{Prob}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$$
- $H_t \triangleq (S_0, R_0, A_0 \ldots S_t, R_t, A_t)$, the **trajectory.**
- $O_t = h(H_t)$ is the **observation**, e.g., image, can be IMU reading, lidar scan …

# Concepts (2) [1]

- Policy function $\pi(\cdot)$:

  - Deterministic policy:

    $$a = \pi(s)$$

  - Stochastic policy:

    $$\pi(a|s) = \text{Prob}(A_t = a|S_t = s)$$

- The return $G_t$:

  $$G_t \triangleq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

- The discount factor:

  $$\gamma \in [0,1]$$

Definition

- (State-)value function under policy $\pi$, $v_\pi(s)$:
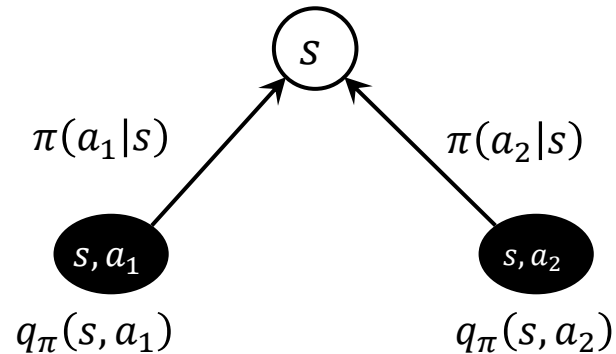
  $$v_\pi(s) \triangleq E_\pi(G_t|S_t = s),$$

- The action-value function:

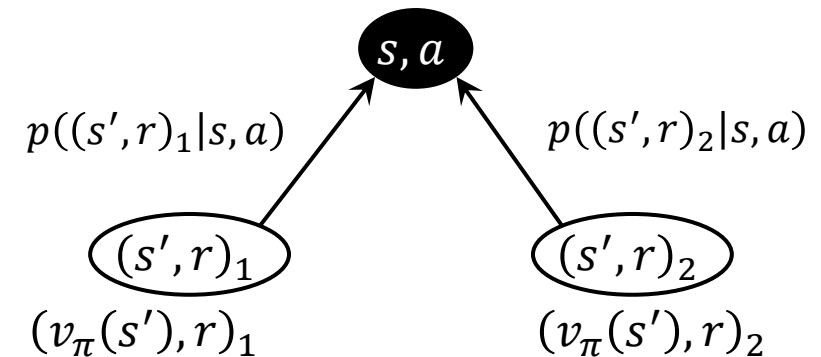  $$q_\pi(s,a) \triangleq E_\pi(G_t|S_t = s, A_t = a),$$

# Bellman Equation

- For state-value function:

$$v_\pi(s) = E_\pi(G_t | S_t = s)$$
$$= \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)$$

- For action-value function:

$$q_\pi(s, a) = E_\pi(G_t | S_t = s, A_t = a) = E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | s, a]$$
$$= \sum_{(s', r) \in \mathcal{S} \times \mathcal{R}} p(s', r | s, a)[r + \gamma v_\pi(s')]$$



$s$

$\pi(a_1|s)$   $\pi(a_2|s)$

$s, a_1$   $s, a_2$

$q_\pi(s, a_1)$   $q_\pi(s, a_2)$

$s, a$

$p((s', r)_1 | s, a)$   $p((s', r)_2 | s, a)$

$(s', r)_1$   $(s', r)_2$

$(v_\pi(s'), r)_1$   $(v_\pi(s'), r)_2$

# Bellman Equation

- For state-value function:

$$v_\pi(s) = E_\pi(G_t|S_t = s)$$

$$= \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)$$

$$= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{(s',r) \in \mathcal{S} \times \mathcal{R}} p(s', r|s, a)[r + \gamma v_\pi(s')]$$

- For action-value function:

$$q_\pi(s, a) = E_\pi(G_t|S_t = s, A_t = a) = E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1})|s, a]$$

$$= \sum_{(s',r) \in \mathcal{S} \times \mathcal{R}} p(s', r|s, a)[r + \gamma v_\pi(s')]$$

$$= \sum_{(s',r) \in \mathcal{S} \times \mathcal{R}} p(s', r|s, a) \left[ r + \gamma \sum_{a' \in \mathcal{A}(s')} \pi(a'|s') q_\pi(s', a') \right]$$

# Optimal Value Functions & BOE

| Definition |
|---|
| ▪ $\pi > \pi' \Rightarrow v_\pi(s) > v_{\pi'}(s), \textcolor{red}{\forall s}$ |
| ▪ The optimal state-value function $v_*(s)$: $$v_*(s) = \max_\pi v_\pi(s)$$ |
| ▪ The optimal action-value function $q_*(s, a)$: $$q_*(s) = \max_\pi q_\pi(s, a)$$ |
| ▪ For any optimal $\pi_*$, all $s \in \mathcal{S}$, all $a \in \mathcal{A}(s)$: $$v_*(s) = \max_{a \in \mathcal{A}(s)} q_*(s, a)$$ $$q_*(s, a) = \sum_{(s', r) \in \mathcal{S} \times \mathcal{R}} p(s', r | s, a)[r + \gamma v_*(s)]$$ |

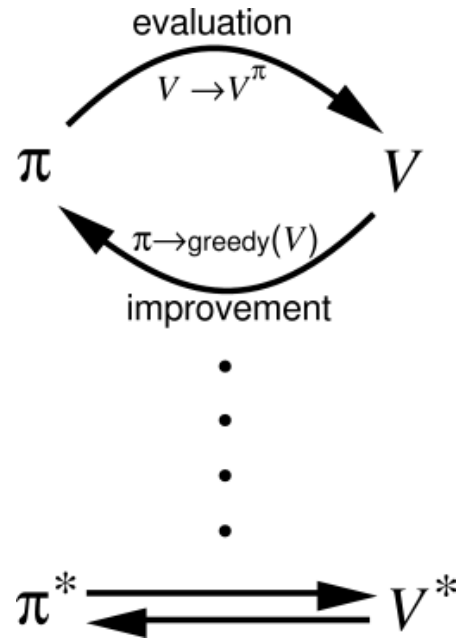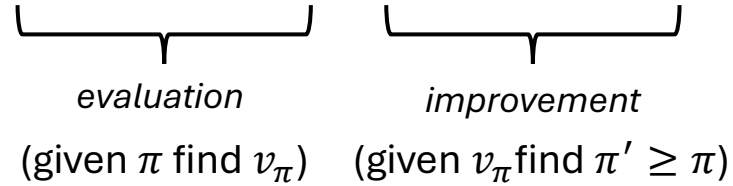| Theorem |
|---|
| For any MDP: |
| ▪ $\exists \pi_*, \pi_* \geq \pi, \forall \pi$ |
| ▪ $v_{\pi_*}(s) = v_*(s), \forall \pi_*$ |
| ▪ $q_{\pi_*}(s, a) = q_\pi(s, a), \forall \pi_*$ |

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount ...
- Evaluation, Iteration, Improvement, Value Iteration ...
- Monte Carlo, Off-policy
- Temporal Difference, Q-learning, Sarsa
- Function Approximation
- Policy Gradient Methods

# Solving the MDP

- *Policy iteration: from some $\pi \to$ evaluate $\pi \to$ improve $\pi$, repeat until $\pi \approx \pi^*$*

evaluation    improvement

(given $\pi$ find $v_\pi$)    (given $v_\pi$ find $\pi' \geq \pi$)

evaluation

$V \to V^\pi$

$\pi$    $V$

$\pi \to$ greedy$(V)$

improvement

$\pi^* \rightleftarrows V^*$

- *Value iteration:* a direct approach that achieves faster convergence.

# Solving the MDP

**Policy *Evaluation*:**
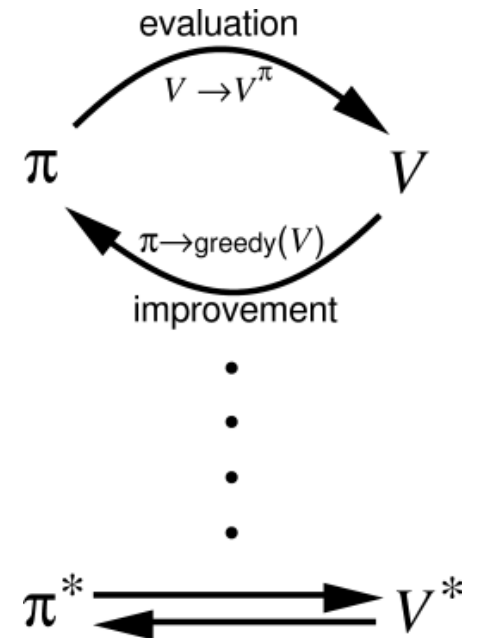
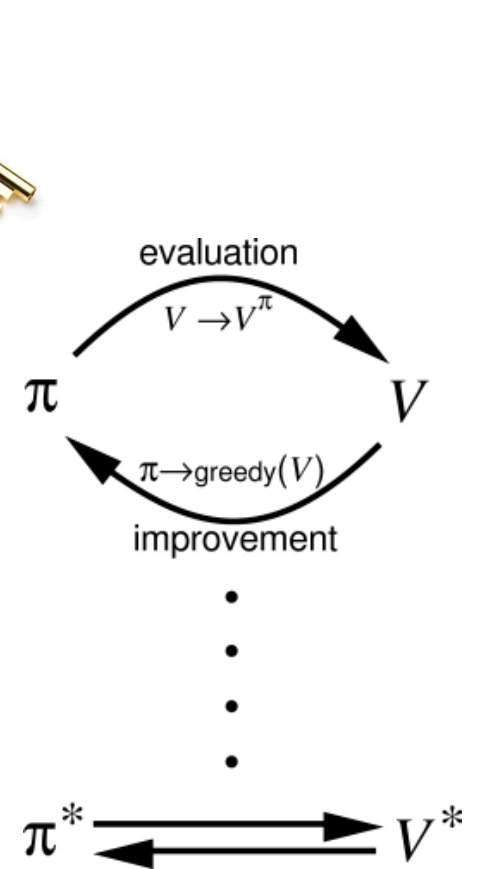Given a policy $\pi(a|s)$
- For $k = 0 \dots K - 1$:

$$\forall s \in \mathcal{S}: V_{k+1}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{(s',r) \in \mathcal{S} \times \mathcal{R}} p(s', r|s, a)[r + \gamma V_k(s')],$$

- $V_k(s) \xrightarrow{K \to \infty} v_\pi(s)$

**Policy *Improvement*:**

Given a value function $v_\pi(s)$:
- $\pi_* = \text{greedy}(v_\pi(s))$



**Policy *Iteration***

# Solving the MDP

**Policy *Evaluation*:**

Given a policy $\pi(a|s)$

- For $k = 0 \dots K-1$:

$$\forall s \in \mathcal{S}: V_{k+1}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{(s',r) \in \mathcal{S} \times \mathcal{R}} \textcolor{red}{p(s',r|s,a)[r + \gamma V_k(s')]},$$

- $V_k(s) \xrightarrow{K \to \infty} v_\pi(s)$

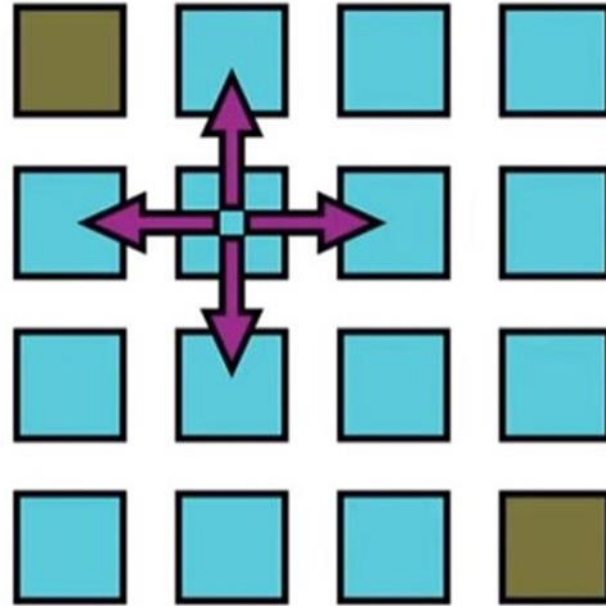**Policy *Improvement*:**

Given a value function $v_\pi(s)$:

- $\pi_* = \text{greedy}(v_\pi(s))$



evaluation

$V \to V^\pi$

$\pi \qquad V$

$\pi \to \text{greedy}(V)$

improvement

$\pi^* \rightleftarrows V^*$

**Policy *Iteration***

# Policy Iteration

Compute $v_\pi(s)$ or $q_\pi(s, a)$ for a given $\pi$.



$R_t = -1$

$\pi(a|s) = 0.25$

# Policy Iteration
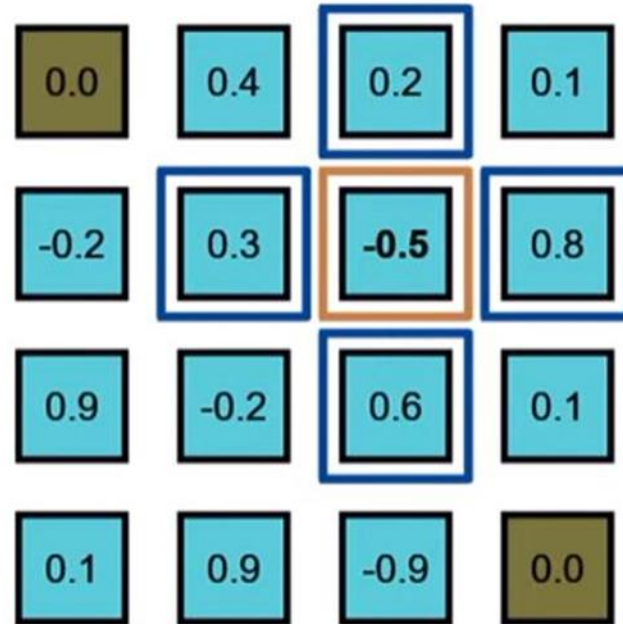
Compute $v_\pi(s)$ or $q_\pi(s, a)$ for a given $\pi$.

| | | | |
|---|---|---|---|
| 0.0 | 0.4 | 0.2 | 0.1 |
| -0.2 | 0.3 | -0.1 | 0.8 |
| 0.9 | -0.2 | 0.6 | 0.1 |
| 0.1 | 0.9 | -0.9 | 0.0 |

$$R_t = -1$$

$$\pi(a|s) = 0.25$$

# Policy Iteration

Compute $v_\pi(s)$ or $q_\pi(s, a)$ for a given $\pi$.



$$R_t = -1$$

$$\pi(a|s) = 0.25$$

$$V(s) \leftarrow \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{\substack{s' \in \mathcal{S} \\ r \in \mathcal{R}}} p(s', r|s, a)[r + \gamma V(s')]$$

# Policy Iteration

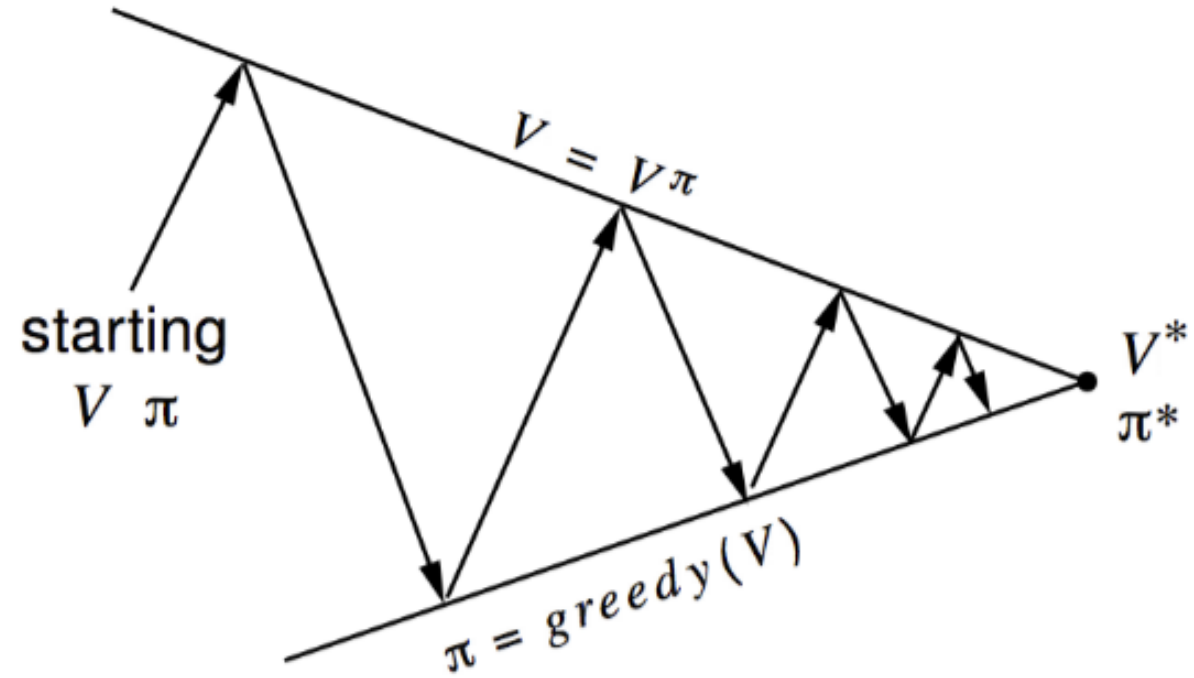Compute $v_\pi(s)$ or $q_\pi(s, a)$ for a given $\pi$.



$$R_t = -1$$

$$\pi(a|s) = 0.25$$

$$V(s) \leftarrow \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{\substack{s' \in \mathcal{S} \\ r \in \mathcal{R}}} p(s', r|s, a)[r + \gamma V(s')]$$

# Policy Iteration



starting $V$ $\pi$

$V = V\pi$

$\pi = greedy(V)$

$V^*$
$\pi^*$

# Value *Iteration*

Find the optimal policy $\pi_*$:

- Given $V_0(s)$:

- Repeat:

- For each $s \in \mathcal{S}$:

- For each $a \in \mathcal{A}(s)$:

- $Q(s,a) \leftarrow \sum_{(s',r)\in\mathcal{S}\times\mathcal{R}} p(s',r|s,a)[r + \gamma V_k(s')]$

- $V_{k+1}(s) \leftarrow \max_{a\in\mathcal{A}(s)} Q(s,a)$

# Tutorial: Tic-Tac-Toe by Value Iteration

**Notes:**

- 'x' goes first w.l.o.g.

- For 3x3 game, neither player can lose if they play optimally:

  → Do not train the AI, play dumb and see that it takes dumb move.

  → Do train the AI, play dumb, and lose to it.

  → Do Train the AI, play smart, and never win over it.

**Notes:**

- $\mathcal{S} = \{1, -1, 0\}^9$

- $R_t = \begin{cases} 1, & \text{if } s_t \text{ win} \\ -1, & \text{if } s_t \text{ loses} \\ 0, & \text{otherwise} \end{cases}$

- $p(s', r | s, a) = \begin{cases} \frac{1}{|\text{legal}(s')|}, & \text{if } (s', r) \text{ is possible} \\ 0, & \text{otherwise} \end{cases}$

- Transition: ...



$s$

$s, a$

$\{s'\}$

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount …

- Evaluation, Iteration, Improvement, Value Iteration …

- Monte Carlo, Off-policy

- Temporal Difference, Q-learning, Sarsa

- Function Approximation

- Policy Gradient Methods

# Monte Carlo Methods

- In real world, most of the time we have imperfect knowledge → estimate.
- Monte Carlo methods are *model-free*

# Monte Carlo Evaluation

- Goal: Given the *data acquired under* $\pi$, estimate $q_\pi$.

- Approach: Express $q_\pi$-estimation problem as $v_\pi$-estimation problem,

    - Define a new problem where:

    $$\bar{S}_t = (S_t, A_t)$$

    $\rightarrow$ Estimating $v(\bar{s})$ is equivalent to estimating $q_\pi(s, a)$.

    - Data = $\left\{ H_m = \left( \bar{s}_0, \bar{s}_1, \dots \bar{s}_{T_m} \right), m = 1 \dots M \right\}$.

    $\rightarrow$ *Markov Reward Process*.

    _____

# Monte Carlo Evaluation

- Goal: Given the *data acquired under* $\pi$, estimate $q_\pi$.

- Approach: Express $q_\pi$-estimation problem as $v_\pi$-estimation problem,

  - Define a new problem where:
  $$\bar{S}_t = (S_t, A_t)$$

  $\rightarrow$ Estimating $v(\bar{s})$ is equivalent to estimating $q_\pi(s, a)$.

  - Data = $\left\{ H_m = (\bar{s}_0, \bar{s}_1, \dots \bar{s}_{T_m}), m = 1 \dots M \right\}$.

  $\rightarrow$ *Markov Reward Process*.

---

- Idea: Use averages to approximate $v_\pi(s) \approx V(s)$:

  - Batch update:

  $$v_\pi(s) = E_\pi(G_t | S_t = s) \approx \frac{1}{C(s)} \sum_{m=1}^{M} \sum_{\tau=0}^{T_m - 1} \mathbb{I}[s_\tau^m = s] \, g_\tau^m \triangleq V(s)$$

# Monte Carlo Evaluation

- Goal: Given the *data acquired under $\pi$*, estimate $q_\pi$.

- Approach: Express $q_\pi$-estimation problem as $v_\pi$-estimation problem,

    - Define a new problem where:

    $$\bar{S}_t = (S_t, A_t)$$

    $\rightarrow$ Estimating $v(\bar{s})$ is equivalent to estimating $q_\pi(s, a)$.

    - Data = $\left\{ H_m = (\bar{s}_0, \bar{s}_1, \dots \bar{s}_{T_m}), m = 1 \dots M \right\}$.

    $\rightarrow$ *Markov Reward Process*.

    _____

- Idea: Use averages to approximate $v_\pi(s) \approx V(s)$:

    - Batch update:

    $$v_\pi(s) = E_\pi(G_t | S_t = s) \approx \frac{1}{C(s)} \sum_{m=1}^{M} \sum_{\tau=0}^{T_m - 1} \mathbb{I}[s_\tau^m = s] \, g_\tau^m \triangleq V(s)$$

    - Iterative update after the $m$-th sample:

    $$V(s_t^m) \leftarrow V(s_t^m) + \frac{1}{C(s_t^m)} \left( g_t^m - V(s_t^m) \right)$$

    - Or simply use a constant step size:

    $$V(s_t^m) \leftarrow V(s_t^m) + \alpha \left( g_t^m - V(s_t^m) \right)$$

$$
\begin{aligned}
\mu_k &= \frac{1}{k} \sum_{j=1}^{k} x_j \\
&= \frac{1}{k} \left( x_k + \sum_{j=1}^{k-1} x_j \right) \\
&= \frac{1}{k} \left( x_k + (k-1)\mu_{k-1} \right) \\
&= \mu_{k-1} + \frac{1}{k} \left( x_k - \mu_{k-1} \right)
\end{aligned}
$$

# MC Control

**Constant-$\alpha$ MC for estimating $\pi \approx \pi*$**

Algorithm inputs:

$$\epsilon \qquad\qquad \alpha \qquad\qquad M$$

Initialize arbitrarily:

$\pi \leftarrow$ some $\epsilon$-soft policy

$Q(s, a) \leftarrow$ some value for $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

For $m = 1, \cdots, M$:

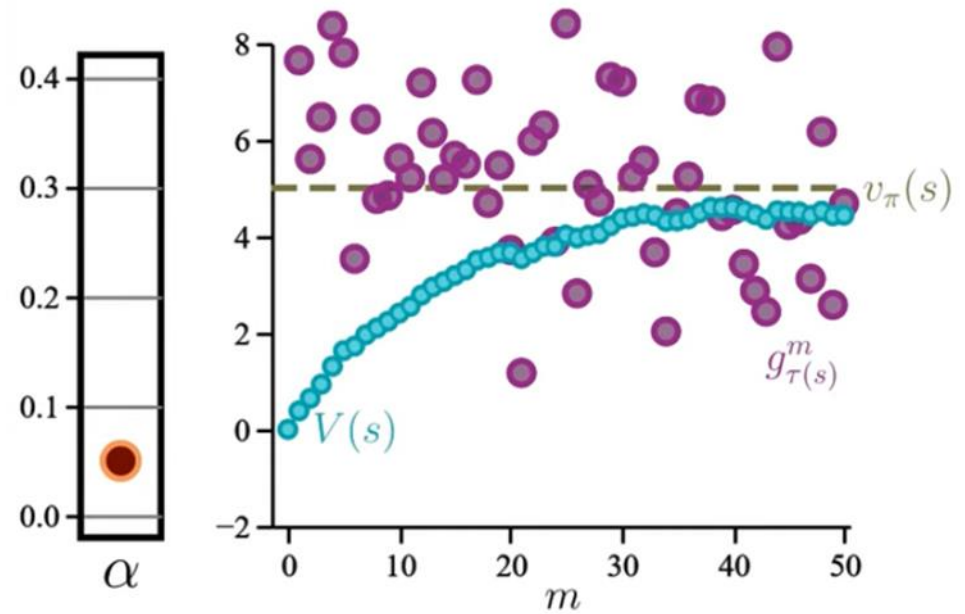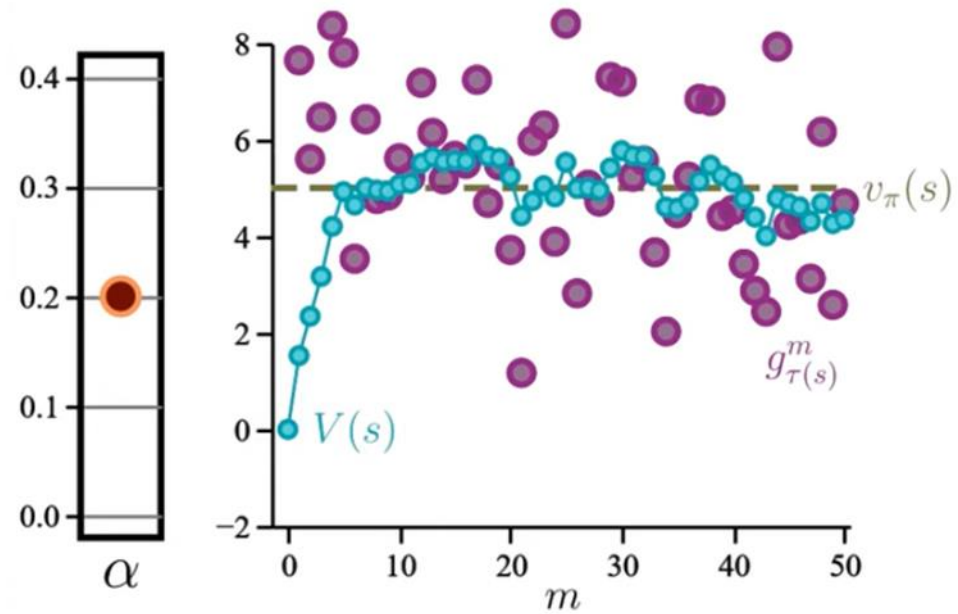Under $\pi$ sample: $s_0^m$, $a_0^m$, $r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

For $t = 0, \cdots, T_m - 1$:

$g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m))$

$\pi \leftarrow \epsilon\text{-greedy}(Q)$
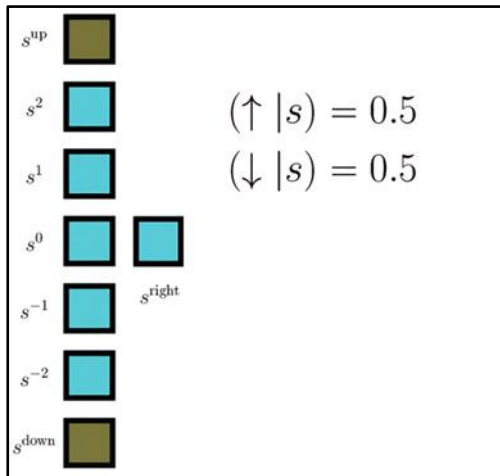
# Monte Carlo Evaluation

# Caveats of MC

- Trajectories have to terminate
- Exploration-Exploitation *dichotomy*:
  - To discover optimal policies, we must **explore** all state-action pairs.
  - To get high returns we must **exploit** known high state-action pairs.

# Caveats of MC

- Trajectories have to terminate
- Exploration-Exploitation dichotomy:
  - To discover optimal policies, we must **explore** all state-action pairs.
  - To get high returns we must **exploit** known high state-action pairs.



With infinite data, $\pi_*$ is always discoverable if the policy is **SOFT**:

$$\pi(a|s) > 0 \qquad \forall s \in \mathcal{S} \quad \forall a \in \mathcal{A}(s)$$

$\epsilon$-**GREEDY POLICY OF** $Q$: With probability $\epsilon$, take an action selected uniformly from $\mathcal{A}(s)$, otherwise take $\mathrm{argmax}_a\, Q(s, a)$.

# Off-policy method

- Goal:

    Estimate $q_\pi(s,a) = E_\pi[G_t | S_t = s, A_t = a]$

- Issue:

    Data exists but collected under $b$

- Remedy:

$$q_\pi(s,a) = E_b\left[\frac{p_\pi(G_t)}{p_b(G_t)} G_t | S_t = s, A_t = a\right]$$

$$\rho = \prod_{\tau=t+1}^{T-1} \frac{\pi(A_\tau | S_\tau)}{b(A_\tau | S_\tau)}$$

$$\pi(a,s) > 0 \Rightarrow b(a,s) > 0$$

**BEHAVIOR POLICY:** Generates the data:

$$b(a|s)$$

**TARGET POLICY:** To be improved/evaluated:

$$\pi(a|s)$$

| ON-POLICY METHODS $b = \pi$ | OFF-POLICY METHODS $b \neq \pi$ |
|---|---|

# Off-policy MC

<u>Constant-$\alpha$ MC for estimating $\pi \approx \pi*$</u>

Algorithm inputs:

$\epsilon$ $\qquad$ $\alpha$ $\qquad$ $M$

Initialize arbitrarily:

$\pi \leftarrow$ some $\epsilon$-soft policy

$Q(s, a) \leftarrow$ some value for $s \in \mathcal{S}, a \in \mathcal{A}(s)$

For $m = 1, \cdots, M$:

Under $\pi$ sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

For $t = 0, \cdots, T_m - 1$:

$g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m))$

$\pi \leftarrow \epsilon\text{-greedy}(Q)$

<u>Off-Policy Constant-$\alpha$ MC for $\pi \approx \pi*$</u>

Algorithm inputs:

$b$ $\qquad$ $\alpha \in (0, 1]$ $\qquad$ $M \in \mathbb{N}$

Initialize arbitrarily:

$\pi \leftarrow$ some policy

$Q(s, a) \leftarrow$ some value for $s \in \mathcal{S}, a \in \mathcal{A}(s)$

For $m = 1, \cdots, M$:

Under $b$ sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

For $t = 0, \cdots, T_m - 1$:

$\rho_t^m \leftarrow \prod_{\tau=t+1}^{T_m-1} \frac{\pi(a_\tau^m|s_\tau^m)}{b(a_\tau^m|s_\tau^m)}$ $\quad$ (or 1 if $t+1 > T_m - 1$)

$g_t^m \leftarrow \rho_t^m(r_{t+1}^m + \gamma r_{t+2}^m + \cdots)$

$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m))$

$\pi(s_t^m) \leftarrow \text{argmax}_a Q(s_t^m, a)$ (ties broken arbitrarily)

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount ...

- Evaluation, Iteration, Improvement, Value Iteration ...

- Monte Carlo, Off-policy

- Temporal Difference, Q-learning, Sarsa

- Function Approximation

- Policy Gradient Methods

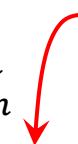# Temporal Difference Learning

- *a priori:*

$$q_\pi(s,a) = E_\pi[\textcolor{red}{G_t}|(S_t, A_t) = (s,a)] = E_\pi[\textcolor{red}{R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1})}|(S_t, A_t) = (s,a)]$$

- *Just read through the maths: $Q(s,a) \approx g_t \approx r_t + \gamma Q(s_{t+1}, a_{t+1})$*

- MC approach:

$$g_t^m = r_{t+1}^m + \gamma r_{t+2}^m \dots \gamma^{T_m-1} r_{T_m}^m$$

target

$$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha\big(g_t^m - Q(s_t^m, a_t^m)\big)$$

# Temporal Difference Learning

- *a priori:*

$$q_\pi(s,a) = E_\pi[G_t|(S_t,A_t) = (s,a)] = E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1},A_{t+1})|(S_t,A_t) = (s,a)]$$

- *Just read through the maths:* $Q(s,a) \approx g_t \approx r_t + \gamma Q(s_{t+1},a_{t+1})$

- MC approach:

$$g_t^m = r_{t+1}^m + \gamma r_{t+2}^m \dots \gamma^{T_m-1} r_{T_m}^m$$

target

$$Q(s_t^m,a_t^m) \leftarrow Q(s_t^m,a_t^m) + \alpha(g_t^m - Q(s_t^m,a_t^m))$$

- 1-step TD approach:

$$\hat{g}_t^m = r_{t+1}^m + \gamma \boxed{Q(s_{t+1},a_{t+1})}$$ bootstrap

$$Q(s_t^m,a_t^m) \leftarrow Q(s_t^m,a_t^m) + \alpha(\hat{g}_t^m - Q(s_t^m,a_t^m))$$

target

# SARSA

<u>**Constant-$\alpha$ MC for estimating $\pi \approx \pi*$**</u>

Algorithm inputs:

$\epsilon$ $\qquad$ $\alpha$ $\qquad$ $M$

Initialize arbitrarily:

$\pi \leftarrow$ some $\epsilon$-soft policy

$Q(s, a) \leftarrow$ some value for $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

For $m = 1, \cdots, M$:

$\quad$ Under $\pi$ sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

$\quad$ For $t = 0, \cdots, T_m - 1$:

$\quad\quad g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

$\quad\quad Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m))$

$\quad\quad \pi \leftarrow \epsilon\text{-greedy}(Q)$

**On Policy TD Control: n-step SARSA**

Changes:

- Approximating the rewards beyond the $n$-th step with the current value of $Q(s, a)$ (bootstrapping):

$$g_{t:t+n}^m = r_{t+1}^m + \cdots + \gamma^{n-1} r_{t+n}^m + \gamma^n Q(s_{t+n}^m, a_{t+n}^m)$$

$$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_{t:t+n}^m - Q(s_t^m, a_t^m))$$
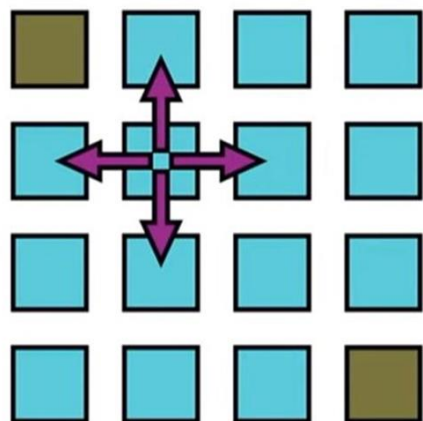
- Updates happen *during* the episode, Interweaving between $(S, A, R)$ tuples, with an n step delay.
- The policy is updated in similar manner with MC

# TD ∋ MC



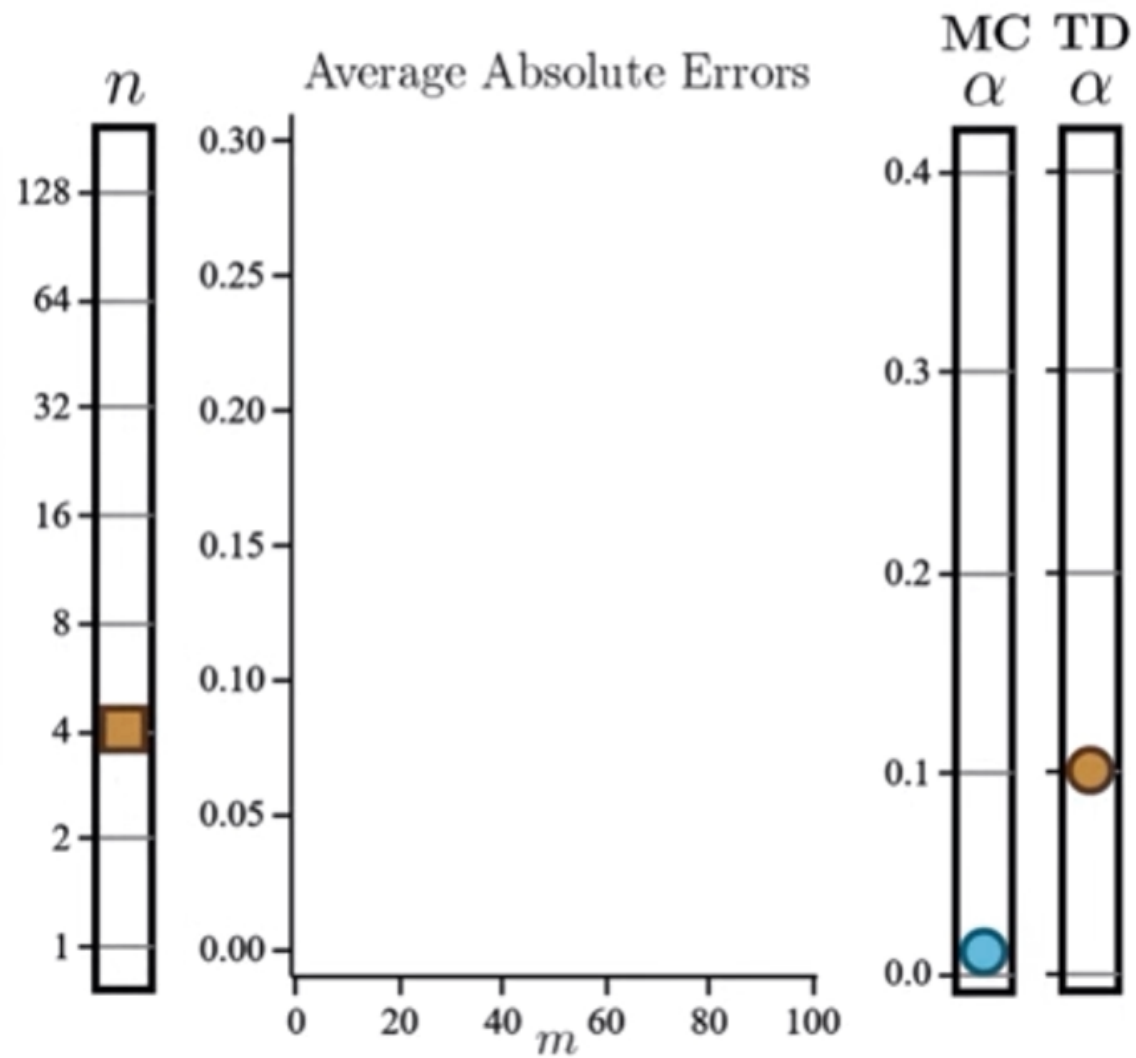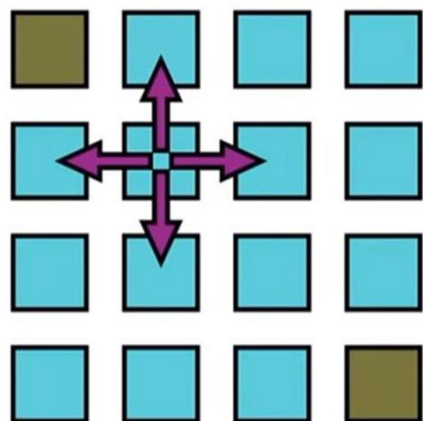Evaluation Example: MC vs TD

# States = 11          # Algo Runs = 200
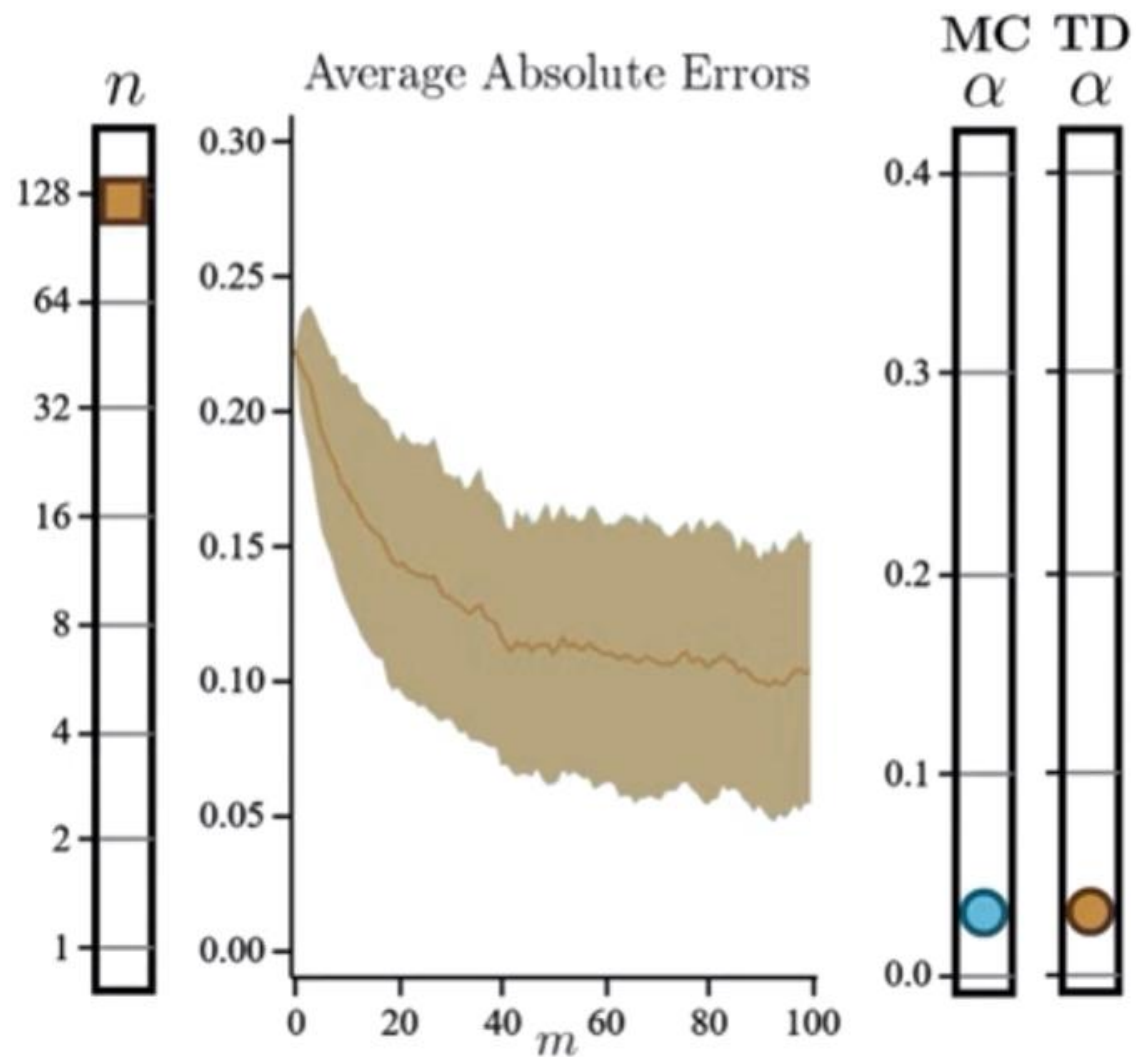
$R_t = -1$

$\pi(a|s) = 0.25$

# TD ∋ MC



$R_t = -1$

$\pi(a|s) = 0.25$

## Evaluation Example: MC vs TD

# States = 11    # Algo Runs = 200

Average Absolute Errors

# Q-learning



**Constant-$\alpha$ MC for estimating $\pi \approx \pi*$**

Algorithm inputs:

$$\epsilon \qquad \alpha \qquad M$$

Initialize arbitrarily:

$\pi \leftarrow$ some $\epsilon$-soft policy

$Q(s,a) \leftarrow$ some value for $s \in \mathcal{S}, a \in \mathcal{A}(s)$

For $m = 1, \cdots, M$:

Under $\pi$ sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

For $t = 0, \cdots, T_m - 1$:

$g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m))$

$\pi \leftarrow \epsilon\text{-greedy}(Q)$

**Q-Learning**

From 1-step **TD** Control, the primary adjustment is to the target:

$$r_{t+1}^m + \gamma Q(s_{t+1}^m, a_{t+1}^m)$$
$$\downarrow$$
$$r_{t+1}^m + \gamma \max_a Q(s_{t+1}^m, a)$$

The max operator means this is **off-policy**.

Under the behavior policy, we are targeting $q_*$.

There's also a change to the update's timing:

1-step **TD:**     update $Q$     update $Q$     update $Q$

$s_0^m, a_0^m, r_1^m, s_1^m, a_1^m, r_2^m, s_2^m, a_2^m, r_3^m, s_3^m, a_3^m, r_4^m \cdots$

1-step **Q:**     update $Q$     update $Q$     update $Q$

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount ...
- Evaluation, Iteration, Improvement, Value Iteration ...
- Monte Carlo, Off-policy
- Temporal Difference, Q-learning, Sarsa
- <span style="color:red">Function Approximation</span>
- Policy Gradient Methods

# Function Approximation

- When $\mathcal{S}$ is continuous $\rightarrow$ never enough data.

- Example:

  - Assume $v_\pi(s) = \hat{v}(s, w), \hat{v}(s + \delta, w) = \hat{v}(s + \delta, w) + \left(\frac{\partial \hat{v}}{\partial s}\right)^\top \delta.$

  - Goal:

  $$\min_{\mathbf{w}} \sum_{s \in \{s_i\}} \|v_\pi(s_i) - \hat{v}(s_i, w)\|^2$$

  - Update rule:

  $$w \leftarrow w + \alpha [G_i - \hat{v}(s_i, w)] \nabla_w \hat{v}(s_i, w)$$

  $$\nabla_w \hat{v}(s_i, w) = \frac{\partial \hat{v}}{\partial w}$$

  - DRL $\rightarrow w$ is param of DNN.

# Function Approximation

- Example:
  - Assume $q_\pi(s) = \hat{q}(s, a, w)$
  - Goal:

$$\min_{\mathbf{w}} \sum_{s \in \{s_i\}} \|q_\pi(s_i, a_i) - \hat{q}(s_i, a_i, w)\|^2$$

  - Update rule:

$$w \leftarrow w + \alpha[G_i - \hat{q}(s_i, a_i, w)] \nabla_w \hat{q}(s_i, a_i, w)$$

$$\nabla_w \hat{q}(s_i, a_i, w) = \frac{\partial \hat{q}}{\partial w}$$

In the context of RL...

- Agent, environment, observations, state, reward, action, value, return, discount ...

- Evaluation, Iteration, Improvement, Value Iteration ...

- Monte Carlo, Off-policy

- Temporal Difference, Q-learning, Sarsa

- Function Approximation

- Policy Gradient Methods

# Policy Gradient Methods

**REINFORCE**

To specify upfront:

- Functional form: $\pi(a|s, \boldsymbol{\theta})$

- Initial $\boldsymbol{\theta}$

- Step size $\alpha$
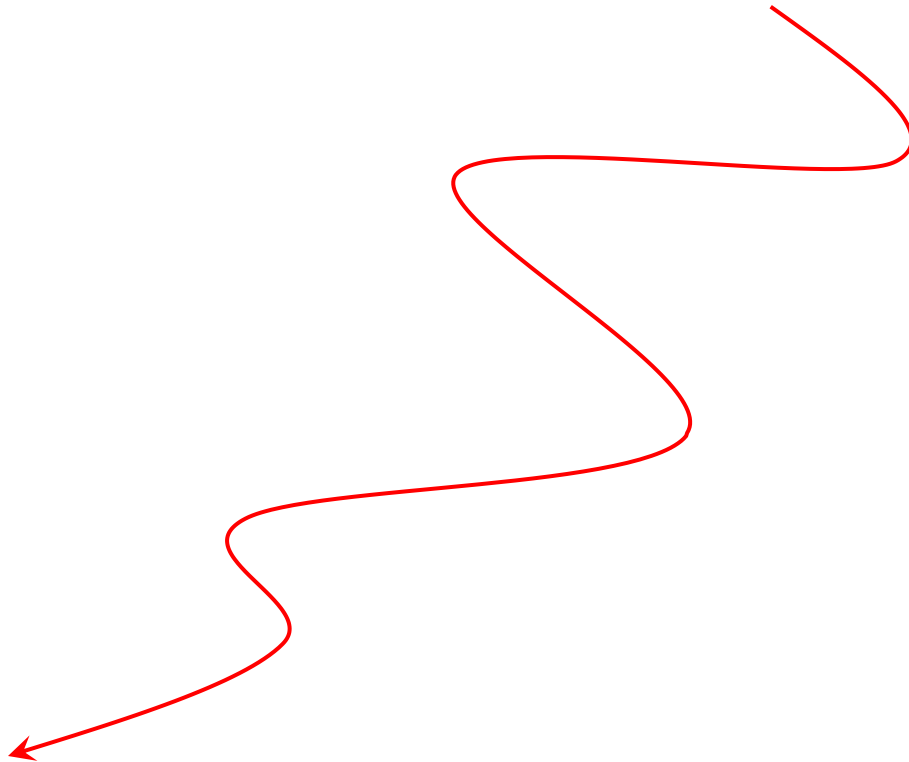
For $m = 1, \cdots, M$:

    Sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

    For $t = 0, \cdots, T_m - 1$:
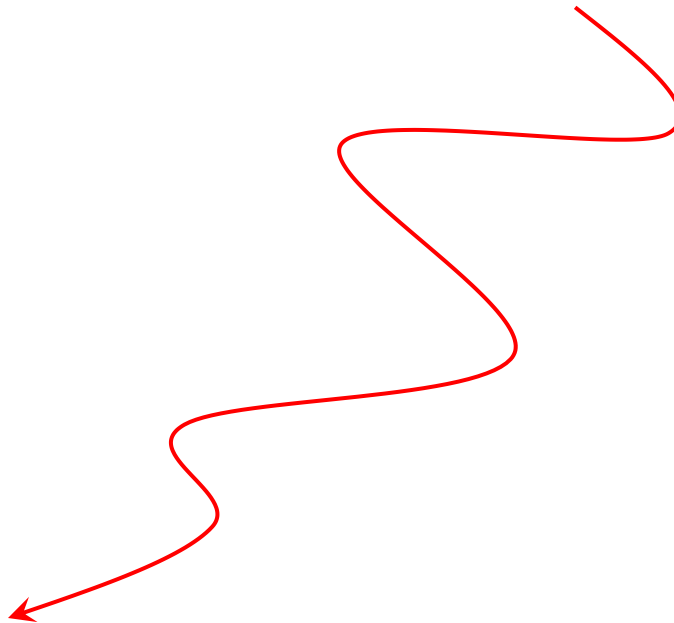
        $g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t g_t^m \nabla \ln \pi(a_t^m | s_t^m, \boldsymbol{\theta})$

- $\nabla_\theta \ln \pi(a_t | s_t, \theta)$ gives the "direction" that increasing $\theta$ will increase $\pi(a_t | s_t, \theta)$.

$$\nabla \ln \pi(a_t^m | s_t^m, \boldsymbol{\theta}) = \frac{\nabla \pi(a_t^m | s_t^m, \boldsymbol{\theta})}{\pi(a_t^m | s_t^m, \boldsymbol{\theta})}$$

# Policy Gradient Methods

**REINFORCE**

To specify upfront:

- Functional form: $\pi(a|s, \boldsymbol{\theta})$

- Initial $\boldsymbol{\theta}$

- Step size $\alpha$

For $m = 1, \cdots, M$:

    Sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

    For $t = 0, \cdots, T_m - 1$:

        $g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$

        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t g_t^m \nabla \ln \pi(a_t^m | s_t^m, \boldsymbol{\theta})$

- $\nabla_\theta \ln\pi(a_t|s_t, \theta)$ gives the "direction" that increasing $\theta$ will increase $\pi(a_t|s_t, \theta)$.

- The increase of $\theta$ is $\sim g_t \nabla_\theta \ln\pi(a_t|s_t, \theta)$

# Policy Gradient Methods

**REINFORCE**

To specify upfront:

- Functional form: $\pi(a|s, \boldsymbol{\theta})$
- Initial $\boldsymbol{\theta}$
- Step size $\alpha$

For $m = 1, \cdots, M$:

Sample: $s_0^m, a_0^m, r_1^m \cdots a_{T_m-1}^m, r_{T_m}^m$

For $t = 0, \cdots, T_m - 1$:

$$g_t^m \leftarrow r_{t+1}^m + \gamma r_{t+2}^m + \cdots$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t g_t^m \nabla \ln \pi(a_t^m | s_t^m, \boldsymbol{\theta})$$

- $\nabla_\theta \ln\pi(a_t|s_t, \theta)$ gives the "direction" that increasing $\theta$ will increase $\pi(a_t|s_t, \theta)$.

- The increase of $\theta$ is $\sim g_t \nabla_\theta \ln\pi(a_t|s_t, \theta)$

$\rightarrow$ the higher the return $g_t$ an action $a_t$ yields, the higher the probability of an action is *increased*.

# Policy Gradient Methods

- **Actor-Critic Methods** combine elements of policy-based methods and value-based methods.

- It introduces an advantage function
$$A(s_i, a_i) = Q(s, a) - V(s)$$
→ provides a measure of how "good" and action is compared with the average action.

- "Actor" gradient:
$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \ln\big(\pi(a_i|s_i, \theta)\big) A(s_i, a_i)$$

- "Critic" gradient:
$$\nabla_w J(w) \approx \frac{1}{N} \sum_i \nabla_w A(s_i, a_i)^2$$

# Tutorial: CartPole

https://www.gymlibrary.dev/environments/classic_control/cart_pole/

# Summary

- Value Evaluation, Policy Iteration, Policy Improvement, Value Iteration...

- MC

- TD, Q-learning, Sarsa

- Function Approximation (Deep RL)

- Policy Gradient Methods

**Model-based RL**

**Model-free RL**

Gradient-free

Off-Policy

On-Policy

$(1)$
$\vdots$
$(\infty) \equiv$

Gradient-based

$$\theta \leftarrow \theta + \alpha \nabla \ln \pi$$

# References

- [Reinforcement Learning: An Introduction, Sutton](#)

- [David Silver RL Lectures](#)

- [Zhao Shiyu RL Lectures](#)

- [OpenAI Introduction to RL](#)