# THE RIGHT-TURN RACER

A foray into Monte Carlo Methods

# PROBLEM STATEMENT



**Figure 5.5:** A couple of right turns for the racetrack task.

- Car has velocity range of 0 to 4, but can only change by -1 to 1
- Car can't stall (velocity vector goes to 0,0)

# SOLUTION

Implement *Monte Carlo method* using an *off-policy* strategy with *weighted importance* sampling.
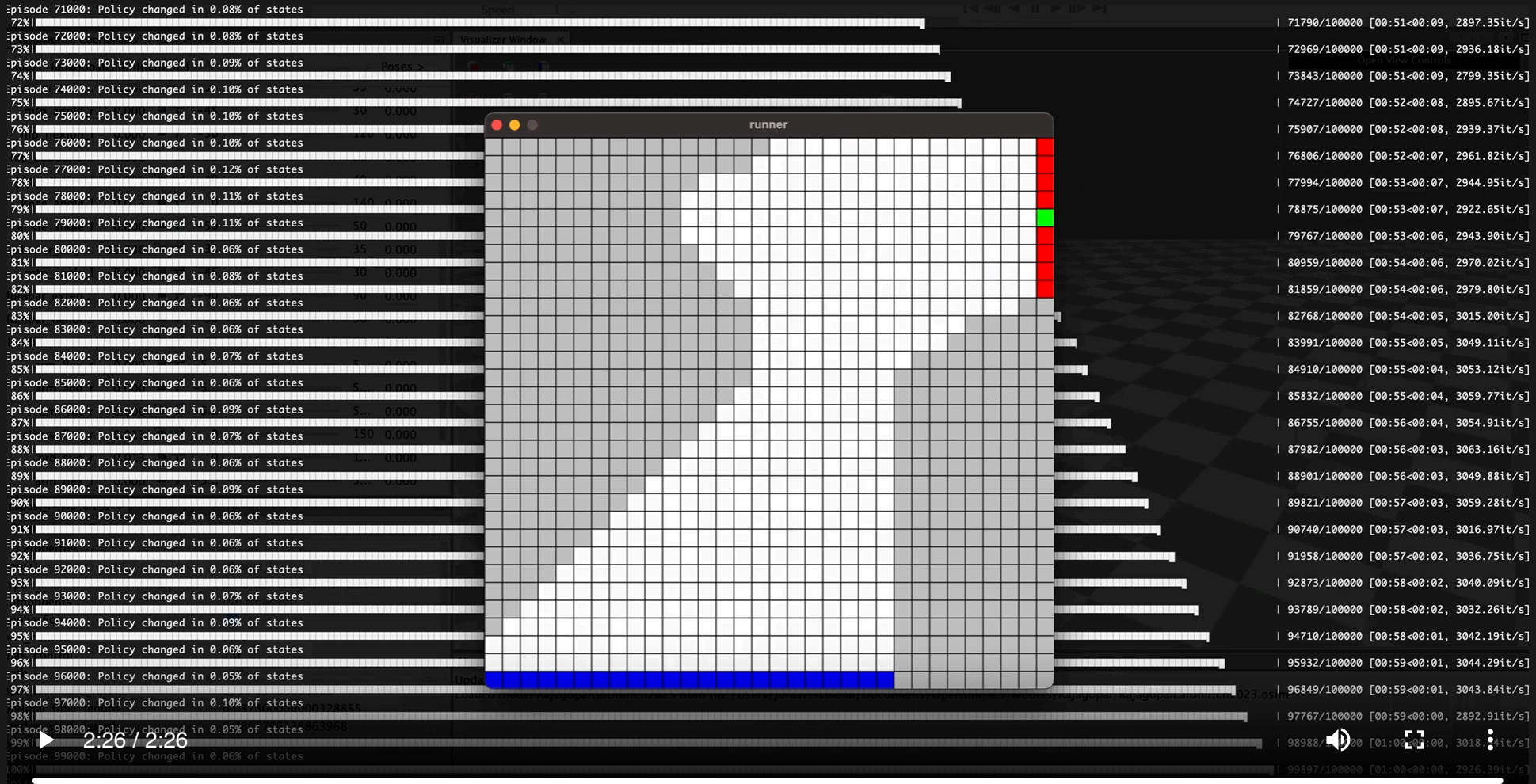
# MONTE CARLO METHOD

- having a model of the world, but no understanding of its *dynamics*
- Must *sample* it rather than raw dog it with math
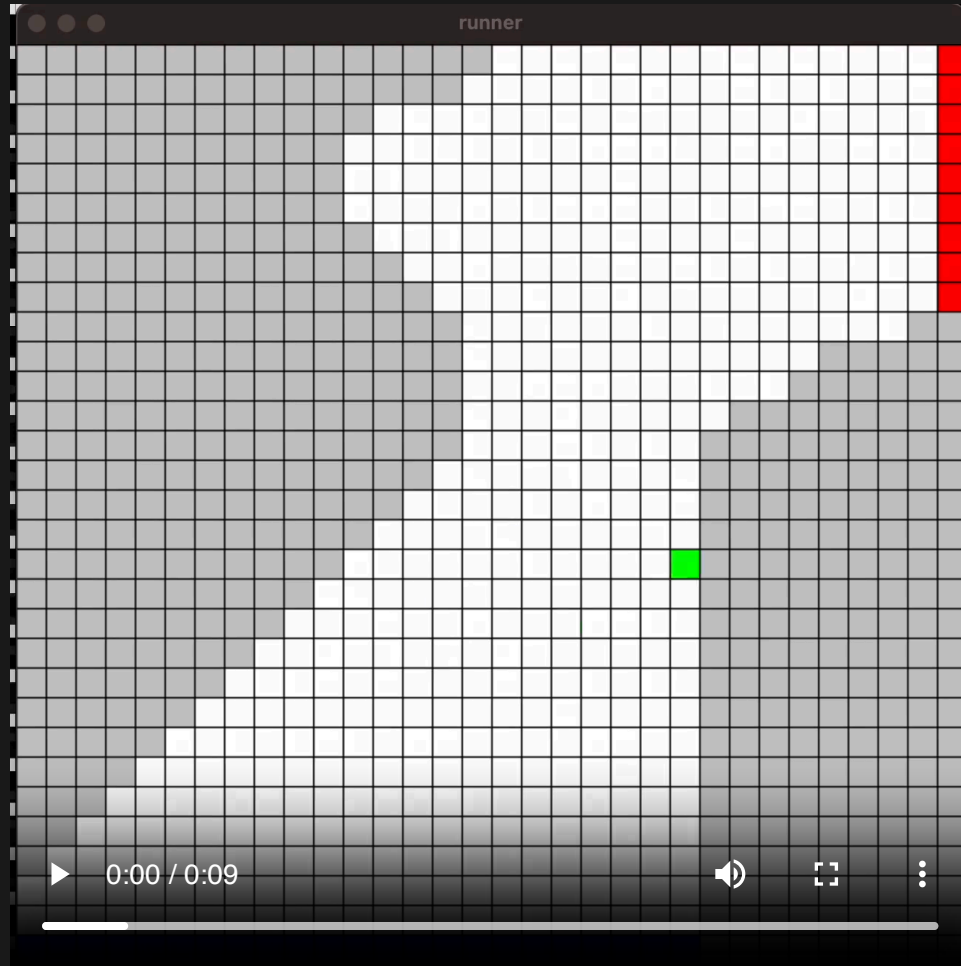
# OFF-POLICY STRATEGY

- the way I choose to *explore* the world will be different than how I choose to *exploit* it
- explore: *behavior policy*
- exploit: *target policy*

## WEIGHTED IMPORTANCE SAMPLING

- my likelihood of seeing the return from this state while exploring is *different* than with my target policy by some *sampling ratio*

But what does this look like while training?

# WHAT PROBLEMS WERE SOLVED FOR ME?

## *Reward Design*

- Decision was to reward -1 per timestep until finish was reached
- Banging into side of track means random restart on start line, NOT END OF EPISODE

# WHAT PROBLEMS DID I HAVE?

*State Design*

Using a simple array to represent the grid made the state space SPARSE, confusing me about the effectiveness of the algorithm

```
  1%|█
pisode 1000: Policy changed in 6.79% of states
  2%|██
pisode 2000: Policy changed in 3.21% of states
  3%|███
pisode 3000: Policy changed in 2.08% of states
  4%|████
pisode 4000: Policy changed in 1.84% of states
  5%|█████
pisode 5000: Policy changed in 1.34% of states
  6%|██████
pisode 6000: Policy changed in 1.46% of states
  7%|███████
pisode 7000: Policy changed in 1.26% of states
  8%|████████
pisode 8000: Policy changed in 1.00% of states
  9%|█████████
pisode 9000: Policy changed in 1.14% of states
 10%|██████████
pisode 10000: Policy changed in 1.09% of states
 11%|███████████
pisode 11000: Policy changed in 0.95% of states
 12%|████████████
pisode 12000: Policy changed in 0.82% of states
 13%|█████████████
pisode 13000: Policy changed in 0.86% of states
 14%|██████████████
pisode 14000: Policy changed in 0.78% of states
 15%|███████████████
pisode 15000: Policy changed in 0.73% of states
 16%|████████████████
pisode 16000: Policy changed in 0.64% of states
 17%|█████████████████
pisode 17000: Policy changed in 0.75% of states
 18%|██████████████████
pisode 18000: Policy changed in 0.61% of states
 19%|███████████████████
pisode 19000: Policy changed in 0.61% of states
 20%|████████████████████
pisode 20000: Policy changed in 0.63% of states
 21%|█████████████████████
pisode 21000: Policy changed in 0.74% of states
 22%|██████████████████████
pisode 22000: Policy changed in 0.64% of states
 23%|███████████████████████
pisode 23000: Policy changed in 0.67% of states
 24%|████████████████████████
pisode 24000: Policy changed in 0.62% of states
 25%|█████████████████████████
pisode 25000: Policy changed in 0.52% of states
 26%|██████████████████████████
pisode 26000: Policy changed in 0.55% of states
 26%|██████████████████████████
```

Visualizer Window ×

Poses >

| | |
|---|---|
| 55 | 0.000 |
| 30 | 0.000 |
| 120 | 0.000 |
| 30 | 0.000 |
| 40 | 0.000 |
| 140 | 0.000 |
| 50 | 0.000 |
| 35 | 0.000 |
| 30 | 0.000 |
| 90 | 0.000 |
| 90 | 0.000 |
| 90 | 0.000 |
| 5... | 0.000 |
| 5... | 0.000 |
| 5... | 0.000 |
| 150 | 0.000 |
| 1... | 0.000 |
| 5... | 0.000 |

Messages × ScriptingShell Window

Updating Model file from 40000 to latest for
Loaded model RajagopalLaiUhlrich2023 from

```
85%|
Episode 85000: Policy changed in 19.78% of states
86%|
Episode 86000: Policy changed in 19.78% of states
87%|
Episode 87000: Policy changed in 19.78% of states
88%|
Episode 88000: Policy changed in 19.78% of states
89%|
Episode 89000: Policy changed in 19.79% of states
90%|
Episode 90000: Policy changed in 19.79% of states
91%|
Episode 91000: Policy changed in 19.80% of states
92%|
Episode 92000: Policy changed in 19.80% of states
93%|
Episode 93000: Policy changed in 19.80% of states
94%|
Episode 94000: Policy changed in 19.81% of states
95%|
Episode 95000: Policy changed in 19.81% of states
96%|
Episode 96000: Policy changed in 19.81% of states
97%|
Episode 97000: Policy changed in 19.83% of states
98%|
Episode 98000: Policy changed in 19.83% of states
99%|
Episode 99000: Policy changed in 19.83% of states
100%|
Episode 100000: Policy changed in 19.84% of states
100%|
```

# WHAT PROBLEMS DID I HAVE?

## *Q-value Initialization*

**Off-policy MC control, for estimating $\pi \approx \pi_*$**

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:
  $\quad Q(s, a) \in \mathbb{R}$ (arbitrarily)
  $\quad C(s, a) \leftarrow 0$
  $\quad \pi(s) \leftarrow \arg\max_a Q(s, a) \quad$ (with ties broken consistently)

Loop forever (for each episode):
  $\quad b \leftarrow$ any soft policy
  $\quad$ Generate an episode using $b$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
  $\quad G \leftarrow 0$
  $\quad W \leftarrow 1$
  $\quad$ Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
    $\quad\quad G \leftarrow \gamma G + R_{t+1}$
    $\quad\quad C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
    $\quad\quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
    $\quad\quad \pi(S_t) \leftarrow \arg\max_a Q(S_t, a) \quad$ (with ties broken consistently)
    $\quad\quad$ If $A_t \neq \pi(S_t)$ then exit inner Loop (proceed to next episode)
    $\quad\quad W \leftarrow W \frac{1}{b(A_t|S_t)}$

# WHAT PROBLEMS DID I HAVE?

*Q-value Initialization*

- The value estimator for the best next action is tied to the reward
- Initializing only positive values made episodes learning useless with a negative only reward

## Q-value Initialization

- Q = [30, 10, 60] at start
  - max is 60, action = 2
- Q = [30, -1, 60] after evaluating final step
  - max is 60, action = 2
    - best action != action taken, ignore rest of episode
- continue until you play russian roulette

# FIN

Code available here

Presentation courtesy of reveal.js