A detailed digital illustration of a desert canyon. The scene is characterized by layered, reddish-brown rock formations and deep, shadowed crevices. Sunlight filters through the upper layers, creating a hazy, golden atmosphere. In the foreground, a lone figure in dark clothing stands on a rocky ledge, looking out over the vast, layered landscape. The overall mood is one of isolation and grandeur.

“One learns from books and example only that certain things can be done. Actual learning requires that you do those things.”

Frank Herbert, *Children of Dune*

CMPUT 365

Introduction to RL

Coursera Reminder

You **should be enrolled in the private session** we created in Coursera for CMPUT 365.

I **cannot** use marks from the public repository for your course marks.

You **need to check, every time**, if you are in the private session and if you are submitting quizzes and assignments to the private section.

The deadlines in the public session **do not align** with the deadlines in Coursera.

If you have any questions or concerns, **talk with the TAs** or email us `cmput365@ualberta.ca`.

Reminders and Notes

- The programming assignment for Control with FA is due today.
- Rich Sutton will give a guest lecture Dec 9th, Monday. Spread the word.
- A note on the final exam:
 - The required reading from the syllabus does not mean that's what will be covered in the final exam. There are some mismatches. Anything we discussed in class is fair game, including Maximization Bias and Double Learning (Section 6.7), and Nonlinear Function Approximation: Artificial Neural Networks (Section 9.7).

Student Perspectives of Teaching (SPOT) Survey

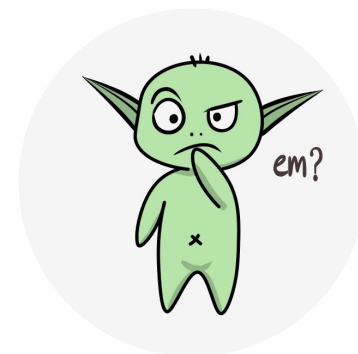


CMPUT 365 LEC A1/A2

Students - FO

<https://go.blueja.io/8c8QQwqLGUqiv9sclG6EDA>

Please, interrupt me at any time!



Last Class: Episodic Semi-gradient One-step Sarsa

Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:

$S, A \leftarrow$ initial state and action of episode (e.g., ε -greedy)

Loop for each step of episode:

Take action A , observe R, S'

If S' is terminal:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

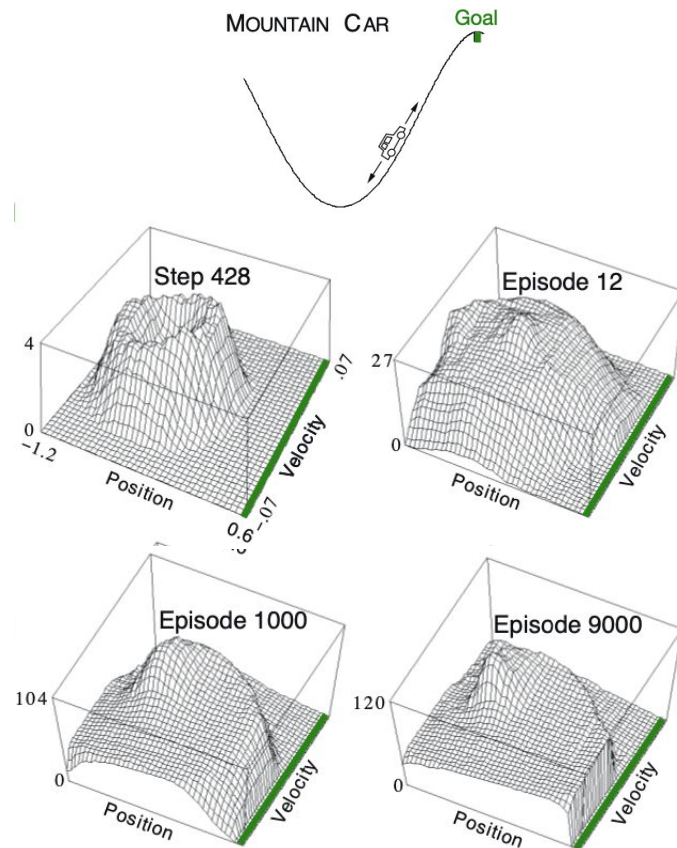
Go to next episode

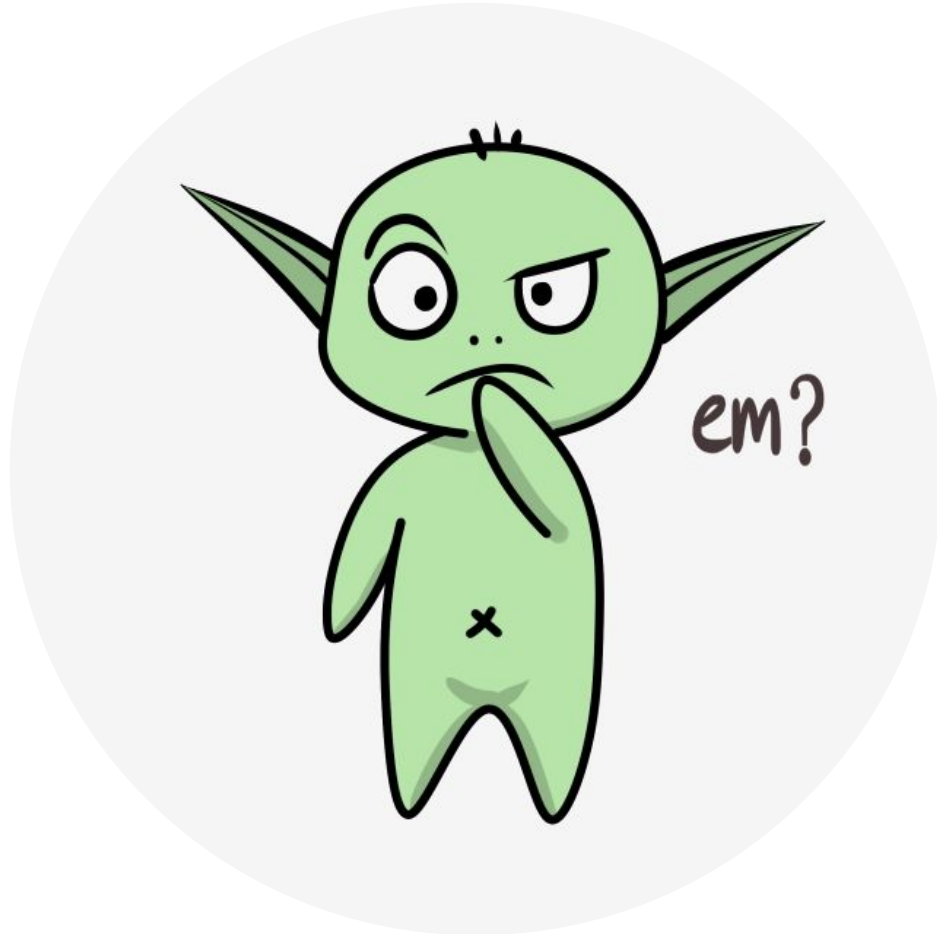
Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

$S \leftarrow S'$

$A \leftarrow A'$





This really works!

State of the Art Control of Atari Games Using Shallow Reinforcement Learning

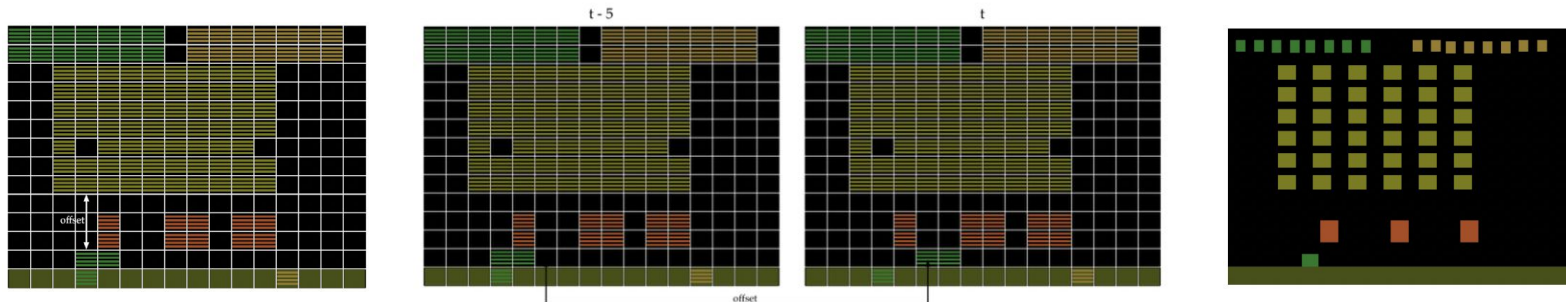
Yitao Liang[†], Marlos C. Machado[‡], Erik Talvitie[†], and Michael Bowling[‡]

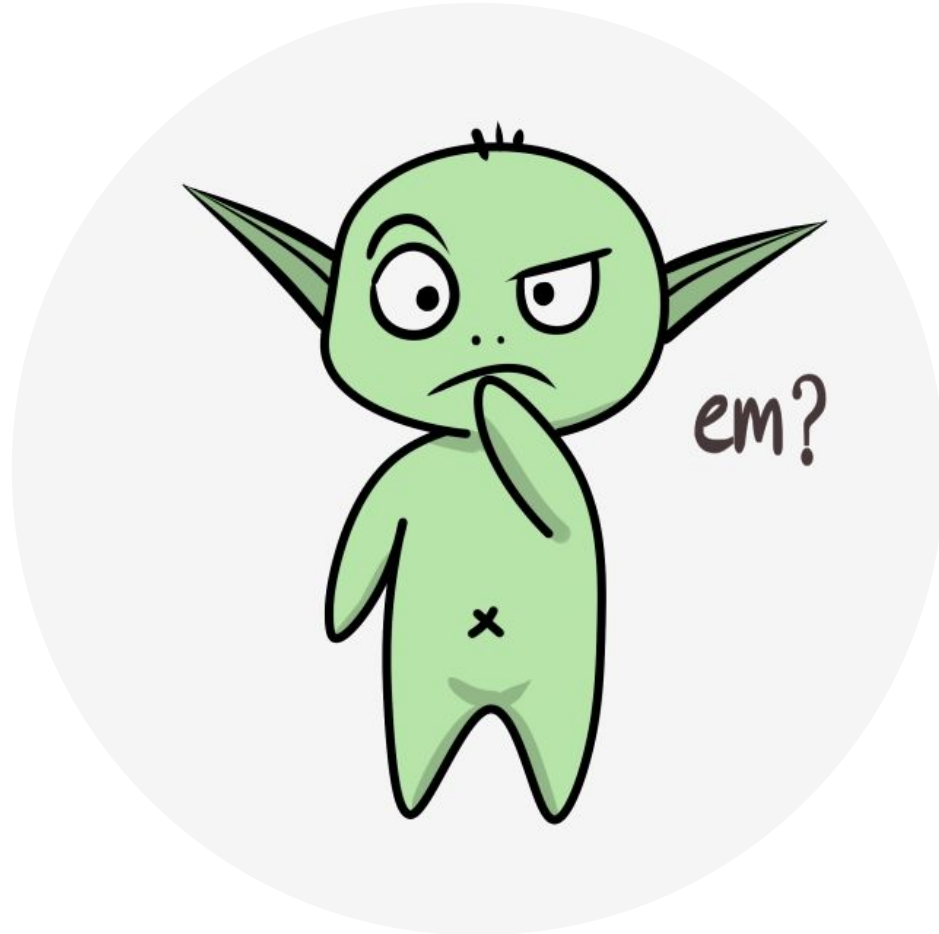
[†]Franklin & Marshall College
Lancaster, PA, USA

[‡]University of Alberta
Edmonton, AB, Canada

{yliang, erik.talvitie}@fandm.edu

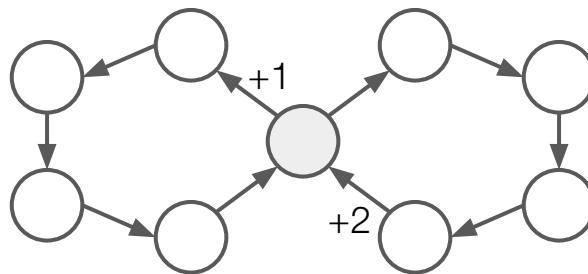
{machado, mbowling}@ualberta.ca





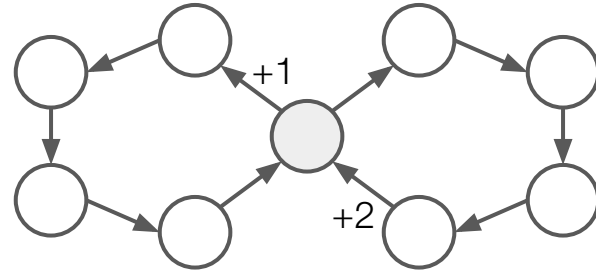
Avg. Reward: A Problem Setting for Continuing Tasks

- Continuing problems without discounting.
 - The agent cares about all rewards equally.



Avg. Reward: A Problem Setting for Continuing Tasks

- Continuing problems without discounting.
 - The agent cares about all rewards equally.



- Quality of a policy is defined by the average rate of reward, $r(\pi)$:

$$\begin{aligned}
 r(\pi) &\doteq \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi] \\
 &= \lim_{t \rightarrow \infty} \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi], \\
 &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) r
 \end{aligned}$$

If the MDP is ergodic: the starting state and any early decision made by the agent can have only a temporary effect; in the long run the expectation of being in a state depends only on the policy and the MDP transition probabilities.

