*"All have their worth,"* said Yavanna,

*"and each contributes to the worth of the others".*

J.R.R. Tolkien, *The Silmarillion*

# CMPUT 365
# Introduction to RL

Marlos C. Machado

# Reminder I

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.

# Reminder II

- Progr. assign. for Coursera's Dynamic Programming module is due Friday.
  Fundamentals of RL: Dynamic Programming – Week 4.

- Monday is a holiday: National Day for Truth and Reconciliation

- Midterm 1 is next Friday.
  Bring your OneCard ID
  No calculators, no cheat sheet

# Last Class: Policy Evaluation (Prediction)

**Iterative Policy Evaluation, for estimating $V \approx v_\pi$**

Input $\pi$, the policy to be evaluated
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$ arbitrarily, for $s \in \mathcal{S}$, and $V(terminal)$ to 0
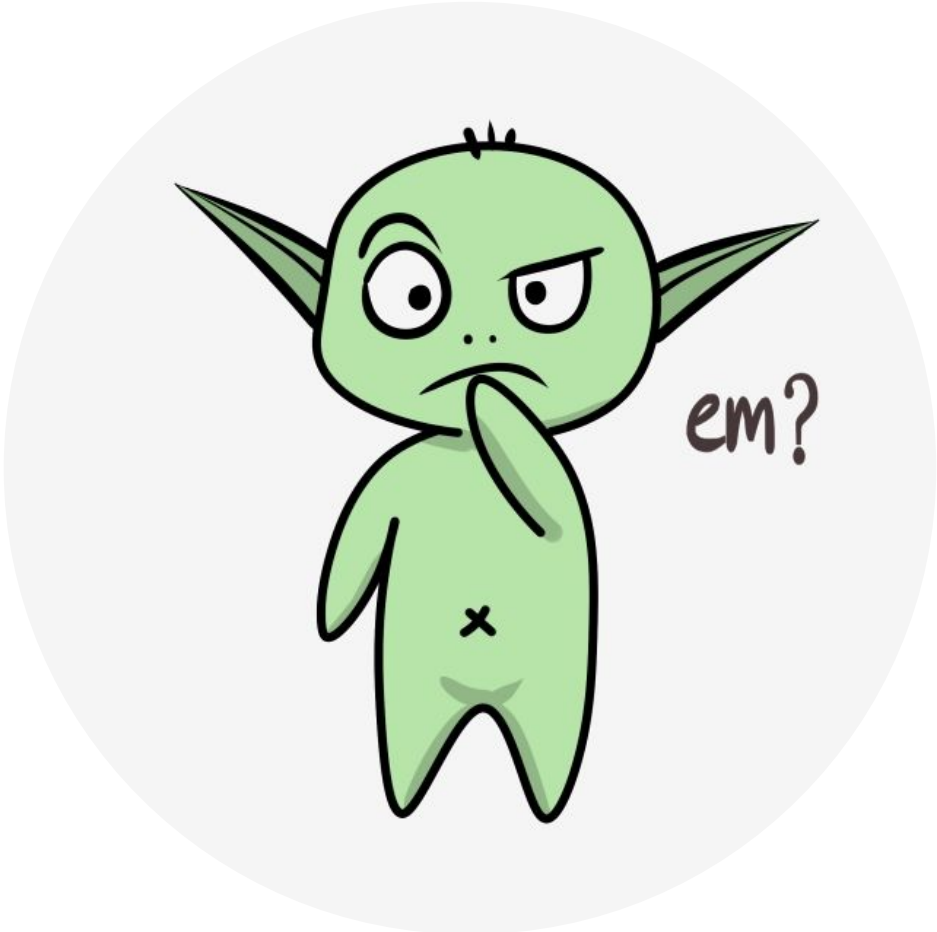
Loop:
    $\Delta \leftarrow 0$
    Loop for each $s \in \mathcal{S}$:
        $v \leftarrow V(s)$
        $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Marlos C. Machado

em?

Marlos C. Machado

# Policy Improvement

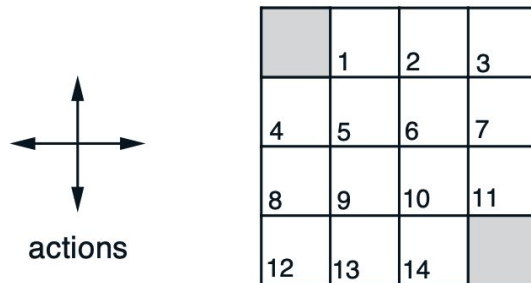*Given a value function for a policy π, how can we get a better policy π'?*

We already know how good policy π is, what if we acted differently now? What if instead of selecting action π(s) we selected action a ≠ π(s), but then we followed π?

We know the value of doing that!

$$q_\pi(s, a) \doteq \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a]$$
$$= \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma v_\pi(s')\right].$$

**If this new action is better, in general this new policy is better overall**

Marlos C. Machado

# Policy Improvement – Intuition



$$R_t = -1$$
on all transitions

$v_k$ for the
random policy

$k = 1$

| 0.0 | -1.0 | -1.0 | -1.0 |
|-----|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

Marlos C. Machado

# Policy Improvement Theorem

That this is true is a special case of a general result called the *policy improvement theorem*. Let $\pi$ and $\pi'$ be any pair of deterministic policies such that, for all $s \in \mathcal{S}$,

$$q_\pi(s, \pi'(s)) \geq v_\pi(s). \tag{4.7}$$

Then the policy $\pi'$ must be as good as, or better than, $\pi$. That is, it must obtain greater or equal expected return from all states $s \in \mathcal{S}$:
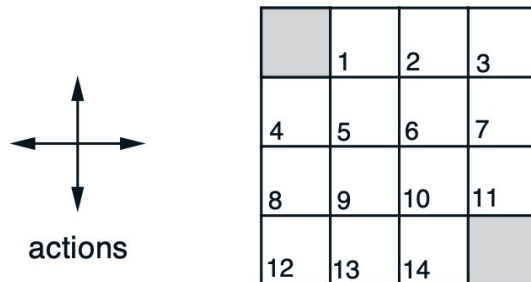
$$v_{\pi'}(s) \geq v_\pi(s). \tag{4.8}$$

# A more aggressive update

Instead of doing it for a particular action in a single state, we can consider changes at *all* states and to *all* possible actions.
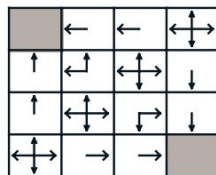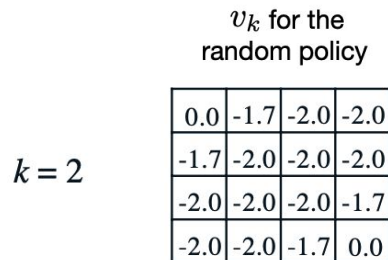
$$
\begin{aligned}
\pi'(s) &\doteq \arg\max_a q_\pi(s, a) \\
&= \arg\max_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\
&= \arg\max_a \sum_{s',r} p(s', r \mid s, a)\Big[r + \gamma v_\pi(s')\Big],
\end{aligned}
$$

This is called *policy improvement*. And eventually it converges to the optimal policy.

# Policy Improvement – Intuition



$$R_t = -1$$
on all transitions

$v_k$ for the
random policy

$k = 2$

| 0.0 | -1.7 | -2.0 | -2.0 |
|------|------|------|------|
| -1.7 | -2.0 | -2.0 | -2.0 |
| -2.0 | -2.0 | -2.0 | -1.7 |
| -2.0 | -2.0 | -1.7 | 0.0 |

actions

Marlos C. Machado

Marlos C. Machado

# Policy Iteration: Interleaving Policy Eval. and Improvement

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$; $V(terminal) \doteq 0$

2. Policy Evaluation
   Loop:
   $\quad$ $\Delta \leftarrow 0$
   $\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad$ $v \leftarrow V(s)$
   $\quad\quad$ $V(s) \leftarrow \sum_{s',r} p(s',r \mid s, \pi(s)) \big[ r + \gamma V(s') \big]$
   $\quad\quad$ $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
   $\quad$ until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

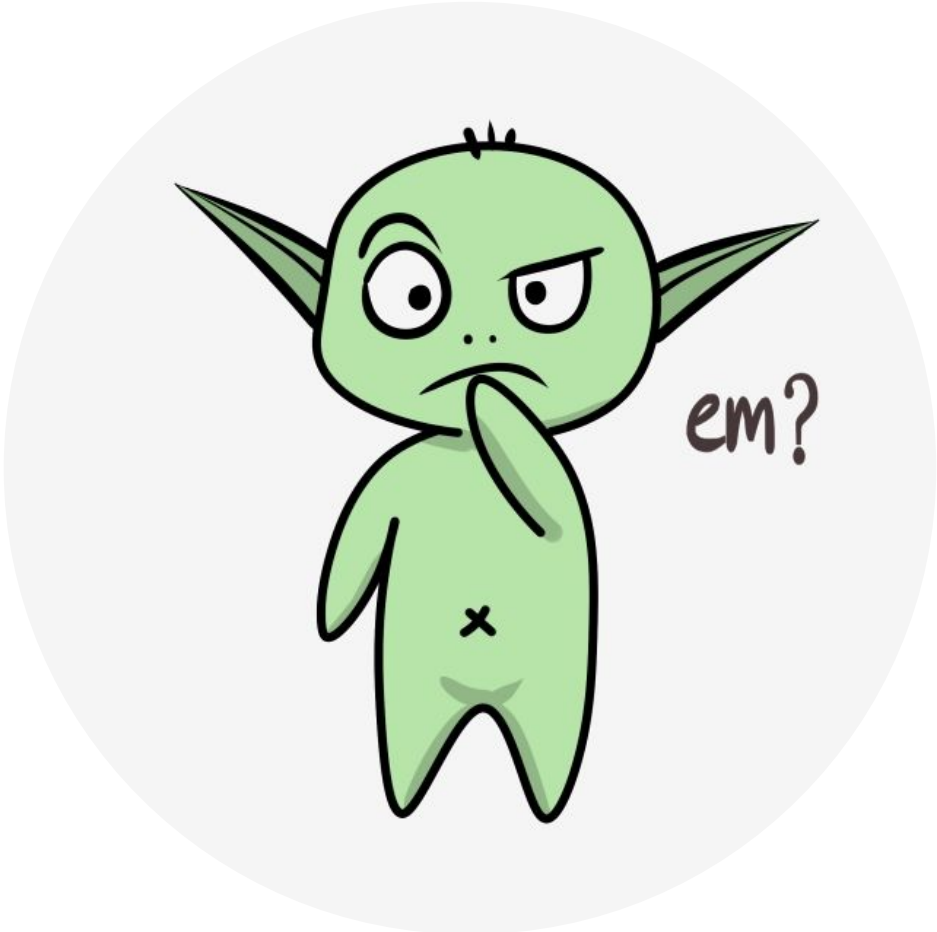3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $\quad$ $old\text{-}action \leftarrow \pi(s)$
   $\quad$ $\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r \mid s, a) \big[ r + \gamma V(s') \big]$
   $\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

# Value Iteration

**Value Iteration, for estimating $\pi \approx \pi_*$**

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
|   $\Delta \leftarrow 0$
|   Loop for each $s \in \mathcal{S}$:
|      $v \leftarrow V(s)$
|      $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
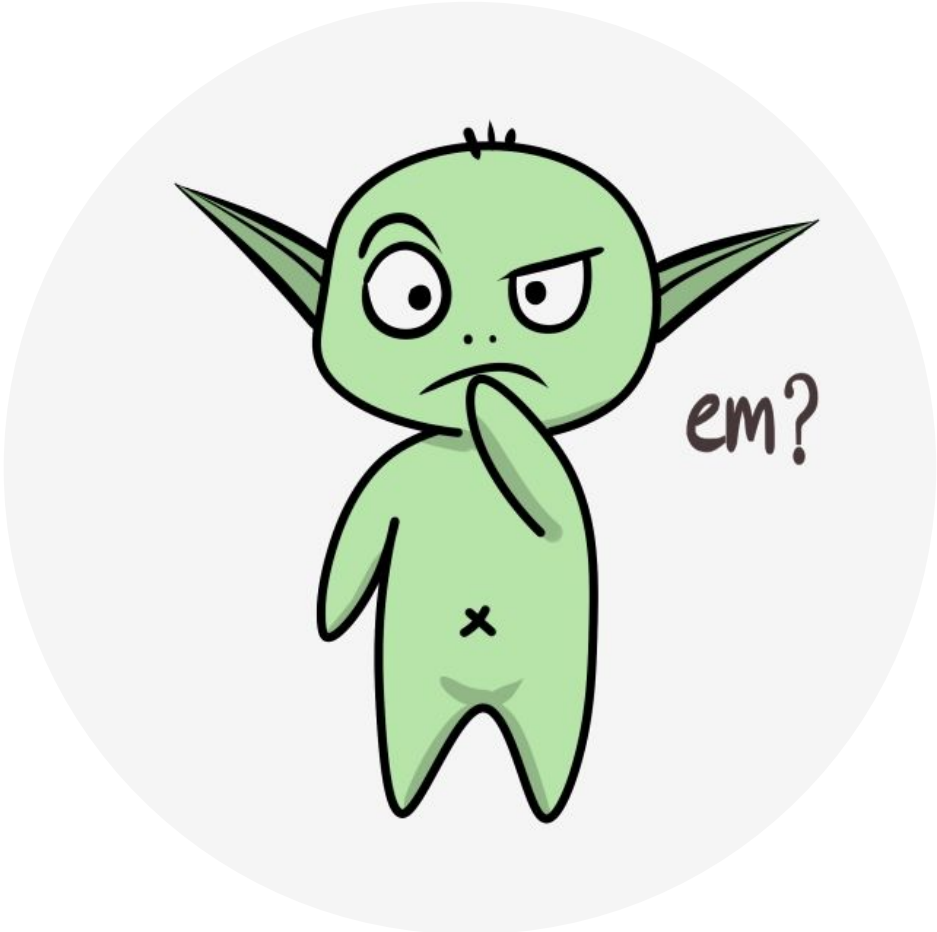|      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
$\pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$

**It doesn't need to be so synchronous**

**We just turned the Bellman optimality equation into an update rule!**

Marlos C. Machado

em?

# Generalized Policy Iteration



evaluation

$V \rightsquigarrow v_\pi$

$\pi$      $V$

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

improvement

$\pi_* \rightleftarrows v_*$

$v = v_\pi$

$v, \pi$

$v_*, \pi_*$

$\pi = \text{greedy}(v)$

em?