

*“The rotten tree-trunk, until the very moment when the storm-blast breaks it in two, has all the appearance of might it ever had.”*

Isaac Asimov, *Foundation*



# **CMPUT 365**

## **Introduction to RL**

# 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.

Some students who are enrolled in Coursera **haven't submitted any quizzes or assignments** in the private session, and that's all I can see.

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

**Please, interrupt me at any time!**



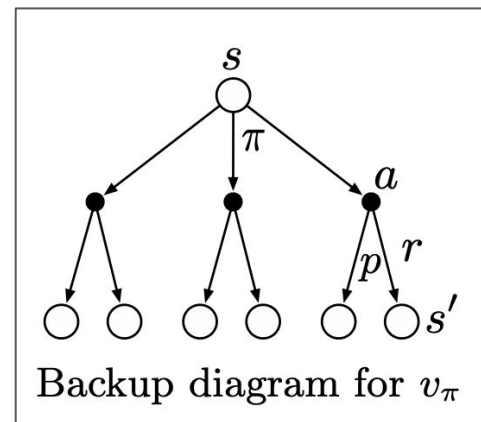
# Value Functions Satisfy Recursive Relationships

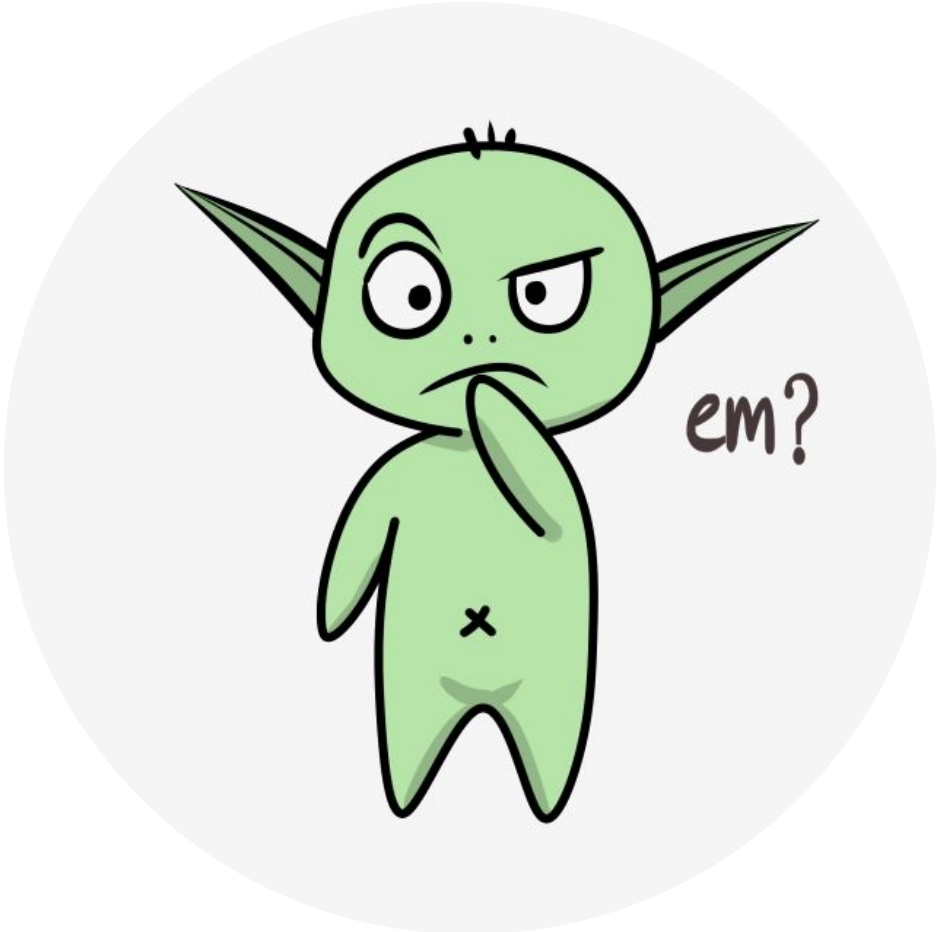
$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

# Value Functions Satisfy Recursive Relationships

$$\begin{aligned}
 v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\
 &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
 &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) \left[ r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) \left[ r + \gamma v_{\pi}(s') \right]
 \end{aligned}$$

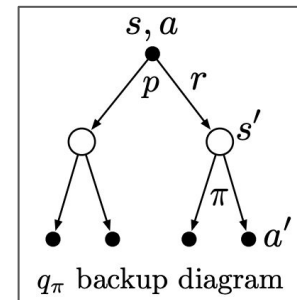
**This is a system of linear equations!**

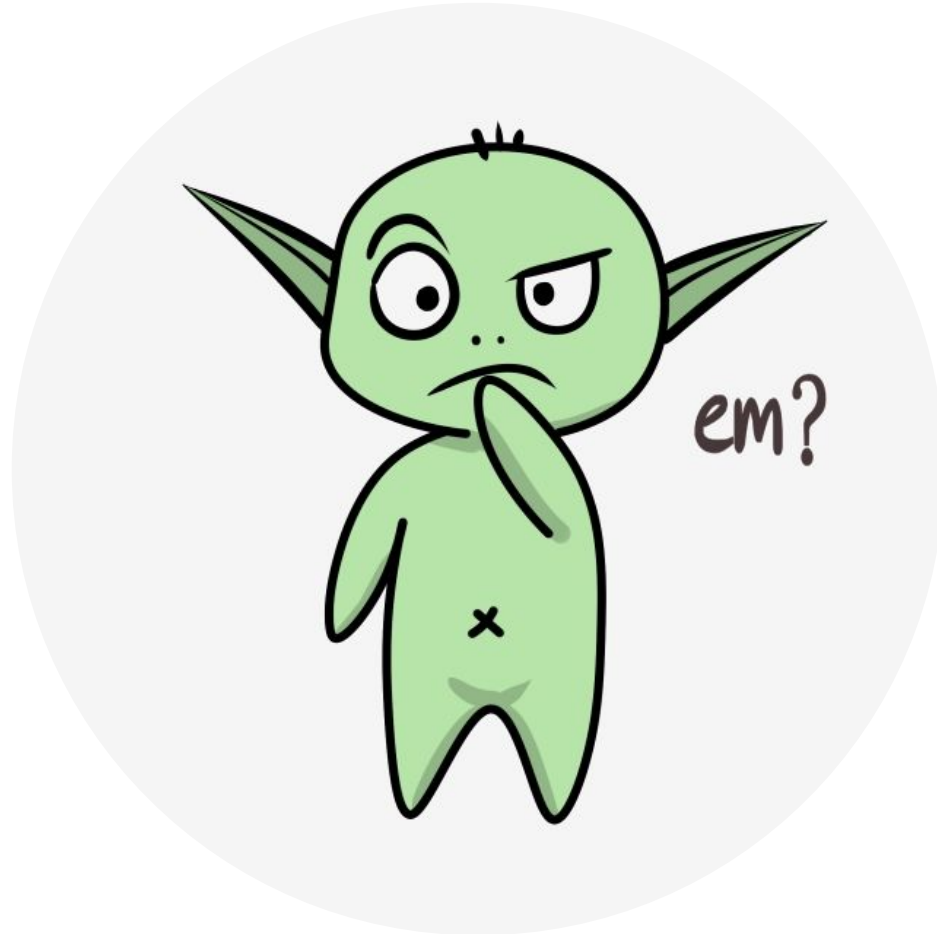




# State-Action Value Functions Satisfy Recursive Relationships

*Exercise 3.17* What is the Bellman equation for action values, that is, for  $q_\pi$ ? It must give the action value  $q_\pi(s, a)$  in terms of the action values,  $q_\pi(s', a')$ , of possible successors to the state–action pair  $(s, a)$ . Hint: The backup diagram to the right corresponds to this equation. Show the sequence of equations analogous to (3.14), but for action values. □







# Optimal Policies and Optimal Value Functions

- Value functions define a partial ordering over policies.
  - $\pi \geq \pi'$  iff  $v_\pi(s) \geq v_{\pi'}(s)$  for all  $s \in \mathcal{S}$ .
  - There is always at least one policy that is better than or equal to all other policies. The *optimal policy*.

$$v_*(s) \doteq \max_{\pi} v_{\pi}(s)$$

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

$$q_*(s, a) \doteq \max_{\pi} q_{\pi}(s, a)$$

# Optimal Policies and Optimal Value Functions

- Because  $v_*$  is the value function for a policy, it must satisfy the self-consistency condition given by the Bellman equation for state values.

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a)$$

# Optimal Policies and Optimal Value Functions

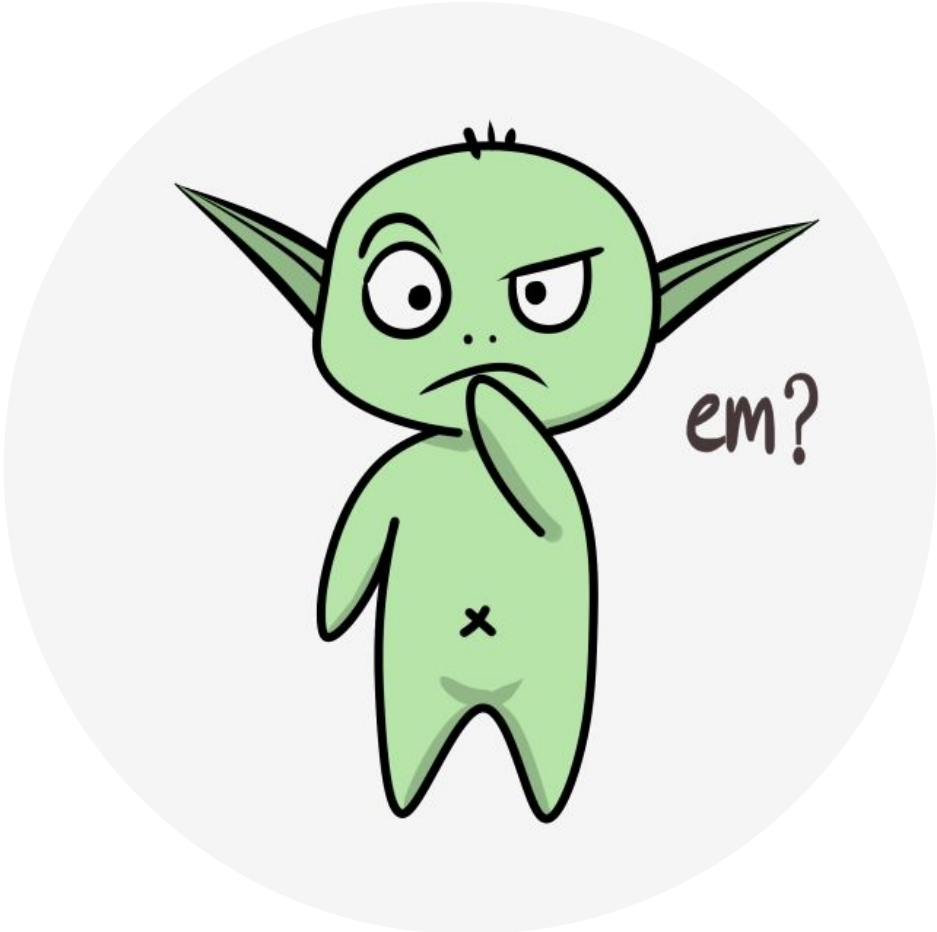
- Because  $v_*$  is the value function for a policy, it must satisfy the self-consistency condition given by the Bellman equation for state values.

$$\begin{aligned}
 v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\
 &= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \\
 &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
 &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\
 &= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')].
 \end{aligned}$$

$$\begin{aligned}
 q_*(s, a) &= \mathbb{E}\left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a\right] \\
 &= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a')\right].
 \end{aligned}$$

# Also...

I have highlighted a couple of exercises during the class, but there are more. The exercises in Chapter 3 of the book are great. I particularly encourage you to look at Exercises 3.25 – 3.29 as well.



## Reinforcement learning is very related to search algorithms

*“Heuristic search methods can be viewed as expanding the right-hand side of the equation below several times, up to some depth, forming a “tree” of possibilities, and then using a heuristic evaluation function to approximate  $v_*$ , at the “leaf” nodes.”*

$$v_*(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')].$$

# Yay! We solved sequential decision-making problems

Except...

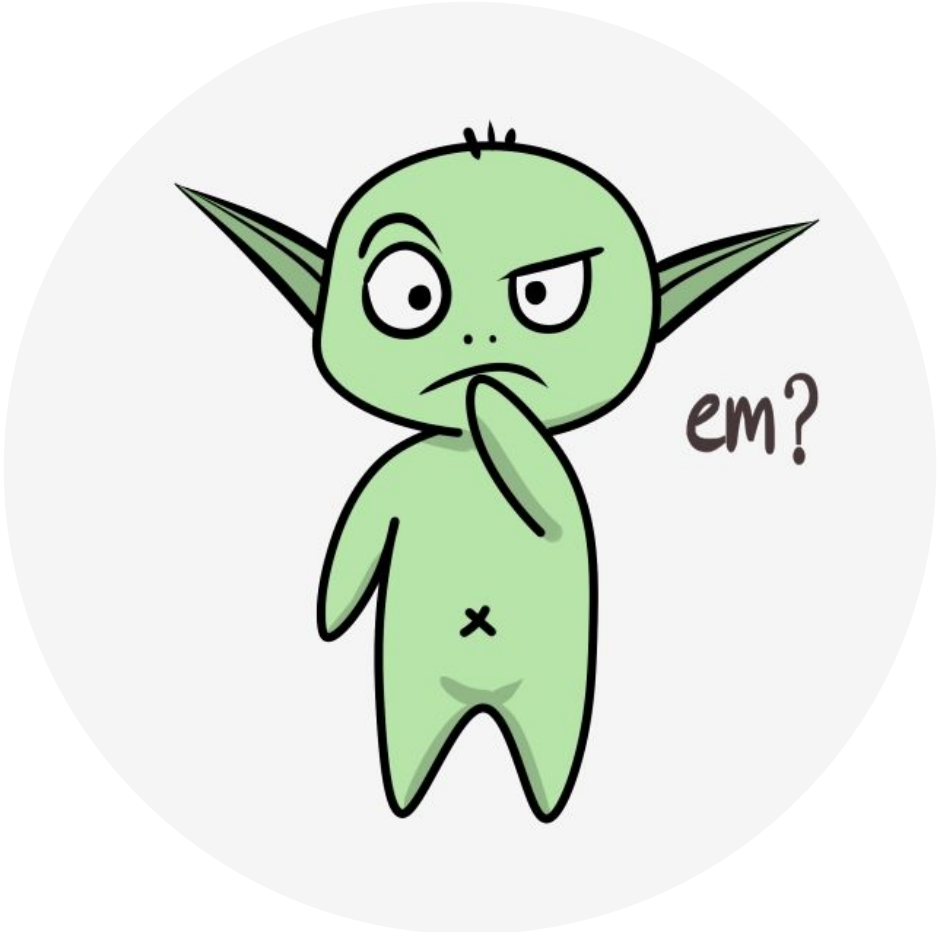
- 1.
- 2.
- 3.

# Yay! We solved sequential decision-making problems

Except...

1. we need to know the dynamics of the environment
2. we have enough computational resources to solve the system of linear eq.
3. the Markov property





# Next class

- What I plan to do:
  - Exercises and Examples
- What I recommend YOU to do for next class:
  - Submit Graded Quiz for Fundamental of RL: Value functions & Bellman equations (Week 3).