

“A beginning is the time for taking the most delicate care that the balances are correct.”

Frank Herbert, *Dune*



CMPUT 365
Introduction to RL

Marlos C. Machado

Class 1 / 35

Plan

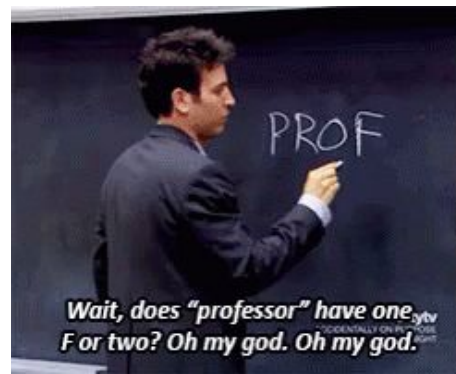
- Introduction
- Course logistics
 - Instruction team
 - Pre-requisites
 - Flipped classroom
 - Textbook
 - Coursera
 - Academic integrity
 - Evaluation
- What is reinforcement learning?

Please, interrupt me at any time!

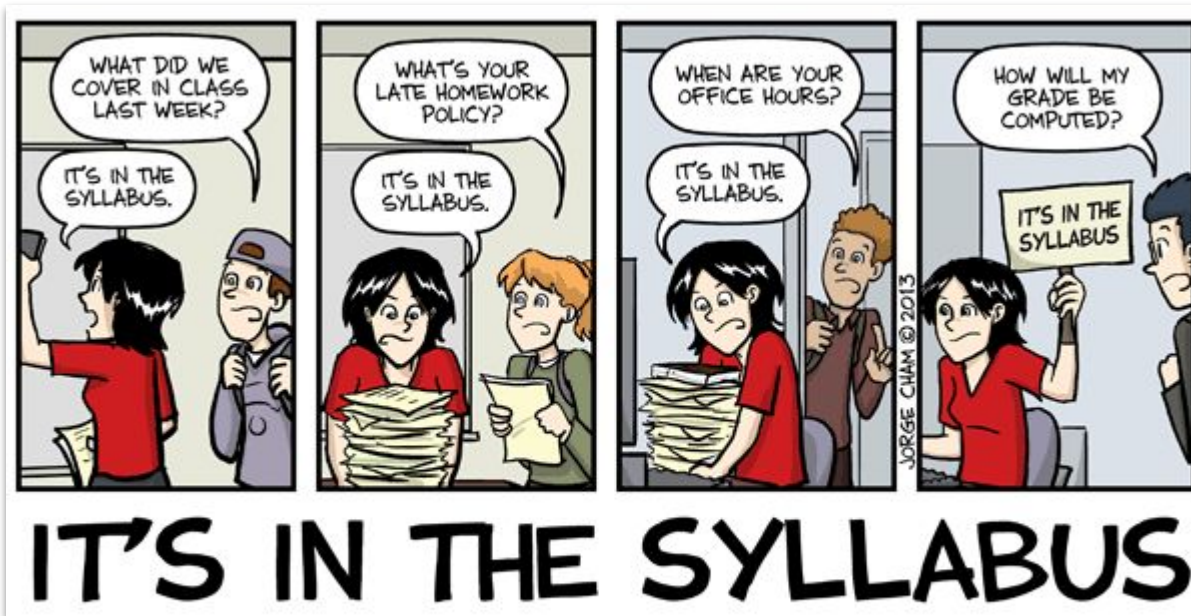


About myself

- Name: Marlos C. Machado
- I was born in Brazil
- I have been living in Edmonton for 10+ years
- I have 2 kids
- Ph.D. working on reinforcement learning
 - Interned at Microsoft Research, IBM Research, and DeepMind
- Worked 4 years at Google Brain and DeepMind
 - Among several other things, we deployed RL to fly balloons in the stratosphere



Course overview and logistics



eClass: [link](#)

Slack: [link](#)

• My website: [link](#)

• Google drive: [link](#)

Start here!

University of Alberta
CMPUT 365: Introduction to Reinforcement Learning
LEC A1
Fall 2024

Instructor: Marlos C. Machado
TAs: Prithvi Nagarajan, Marcos José, Harsh Kotamreddy, Mohamed Mohamed, and Lucas Cruz

Office: ATH 3-08
E-mail: machado@ualberta.ca
Web Page: https://courses.cpsc.ualberta.ca/course/365w_spt14n3924

Office hours: Marlos C. Machado: Thursday 13:00 - 15:00 in ATH 3-08 (Athabasca Hall)
 Prithvi Nagarajan: Monday 11:00 - 13:00
 Lucas Cruz: Tuesday 10:00 - 12:00
 Harsh Kotamreddy: Wednesday 10:00 - 12:00
 Mohamed Mohamed: Thursday 10:00 - 12:00
 Marcos José: Friday 10:00 - 12:00
 The location in which the TAs will hold office hours will be available on eClass.
 Slack and eClass: asynchronously

TA email address: casu365@ualberta.ca
 Do not personally email the TAs. They will only respond via casu365@ualberta.ca.

Lecture room & time: ESB 3-27, MWf 13:00 - 13:20
 Attendance isn't mandatory, although strongly encouraged.

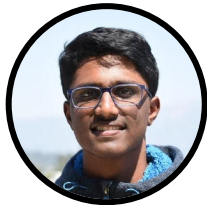
Slack invitation link: We will use Slack as an optional alternative to eClass for communication and question answering. The invitation link will be provided to the students on eClass.

COURSE CONTENT

Course Description: This course provides an introduction to reinforcement learning, which focuses on the study and design of learning agents that interact with a complex, uncertain world to achieve a goal. The course will cover multi-armed bandits, Markov decision processes, reinforcement learning, planning, and function approximation (online supervised learning). The

Key resources

- Syllabus
 - eClass, Slack, my website, Google Drive.
- Teaching assistants



Harshil



Lucas





Marcos



Mohamed



Prabhat

- TA email address: cmput365@ualberta.ca
- My email address: machado@ualberta.ca
-  Slack  invitation link: [link](#)

I want to make this course is a **safe** and **inclusive** environment, for everyone.

It is ok to make mistakes.

We should all strive to be **respectful** to each other.

Office hours

- Slack and eClass: Asynchronous
- Marlos: Thursday 13:00 - 15:00 in ATH 3-08 (Athabasca Hall)
- Prabhat: Mon 11:00 - 13:00 in CSC 2-50
- Lucas: Tue 10:00 - 12:00 in CAB 3-13
- Harshil: Wed 10:00 - 12:00 in CAB 3-13
- Mohamed: Thu 10:00 - 12:00 in CAB 3-13
- Marcos: Fri 10:00 - 12:00 in CAB 3-13

Syllabus [[eClass](#), [Slack](#), [website](#), [Google Drive](#)]

Pre-requisites

- CMPUT 175 or CMPUT 275
- CMPUT 267 or 466, or STAT 265
- Python
- Probability (e.g., expectations of random variables, conditional expectations)
- Calculus (e.g., partial derivatives)
- Linear algebra (e.g., vectors and matrices)

You should either be familiar with these topics or be ready to pick them up quickly as needed by consulting outside resources.

This will be *sort of* a flipped classroom!

- Roughly, you are initially introduced to **new topics outside** the classroom, so we can use the classroom time to explore topics in greater depth
 - A lecture is not necessarily the best use of class time

This will be *sort of* a flipped classroom!

- Roughly, you are initially introduced to **new topics outside** the classroom, so we can use the classroom time to explore topics in greater depth
 - A lecture is not necessarily the best use of class time
- This is about creating meaningful learning opportunities for you, with more personalized interactions – to create **engaged** learning experiences

This will be *sort of* a flipped classroom!

- Roughly, you are initially introduced to **new topics outside** the classroom, so we can use the classroom time to explore topics in greater depth
 - A lecture is not necessarily the best use of class time
- This is about creating meaningful learning opportunities for you, with more personalized interactions – to create **engaged** learning experiences
- I'm not doing this because it is easy, but because I think it is right
 - This is much much more work for me



This will be *sort of* a flipped classroom!

- Roughly, you are initially introduced to **new topics outside** the classroom, so we can use the classroom time to explore topics in greater depth
 - A lecture is not necessarily the best use of class time
- This is about creating meaningful learning opportunities for you, with more personalized interactions – to create **engaged** learning experiences
- I'm not doing this because it is easy, but because I think it is right
 - This is much much more work for me
- This **does not** mean lack of proper guidance, or that you have to teach yourself
- But you do have to become an **active** learner, instead of a passive learner



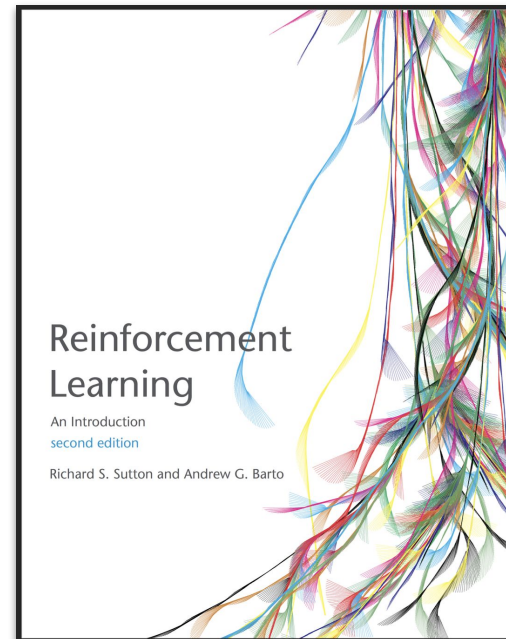
Required textbook

Reinforcement Learning: An Introduction

Richard S. Sutton & Andrew G. Barto

MIT Press. 2nd Edition.

<http://www.incompleteideas.net/book/the-book-2nd.html>



Required textbook

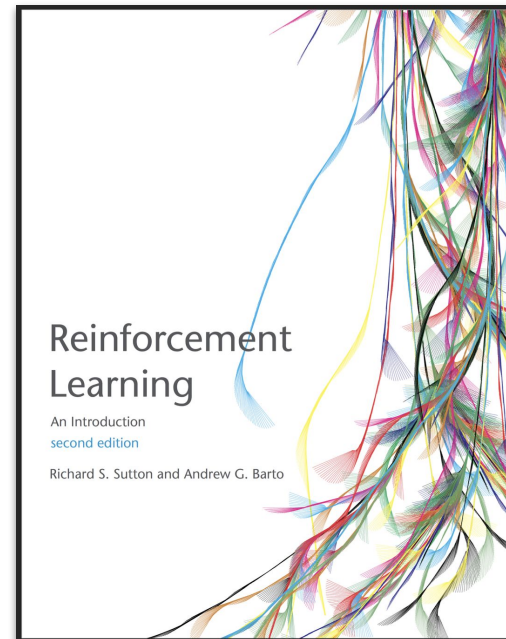
Reinforcement Learning: An Introduction

Richard S. Sutton & Andrew G. Barto

MIT Press. 2nd Edition.

<http://www.incompleteideas.net/book/the-book-2nd.html>

- You will need to read the book!
(This is a sort of a flipped classroom, remember?)
- The book is really good!



GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024
Final exam	30%	December 17, 2024*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024
Final exam	30%	December 17, 2024*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)

Coursera, almost every week (starting Monday): 31.5%

Final exam

30%

December 17, 2024*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)

Coursera, almost every week (starting Monday): 31.5%

Late submissions will not be accepted. There are 11 quizzes and 11 graded assignments. You're expected to do all of them, but s**t happens, so you can miss 2 of each and still get full marks.

Syllabus [[eClass](#), [Slack](#), [website](#), [Google Drive](#)]

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	9 x 1% = 9%	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	9 x 2.5% = 22.5%	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024

Two midterms, summing to 40%. Closed book.

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	9 x 1% = 9%	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	9 x 2.5% = 22.5%	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024

Two midterms, summing to 40%. Closed book.

If you miss the midterm, you can can apply for an excused absence.
If granted, the weight of the missed midterm will be deferred to the final.

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024
Final exam	30%	December 17, 2024*

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)

The final, worth 30%, will be about the whole course.
If you miss the final, you can apply to a deferred final examination.

Final exam	30%	December 17, 2024*
------------	-----	--------------------

GRADE EVALUATION

Assessment	Weight	Date
Practice quizzes (80% pass)	$9 \times 1\% = 9\%$	Day of the 1st class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Assessments (graded quizzes/notebooks on Coursera)	$9 \times 2.5\% = 22.5\%$	Day of the 2nd class on the topic of the week at 23:59:59 (see Course schedule at the end for details)
Midterm 1 exam	20 %	October 4, 2024
Midterm 2 exam	20%	November 1, 2024
Final exam	30%	December 17, 2024*

GRADE EVALUATION

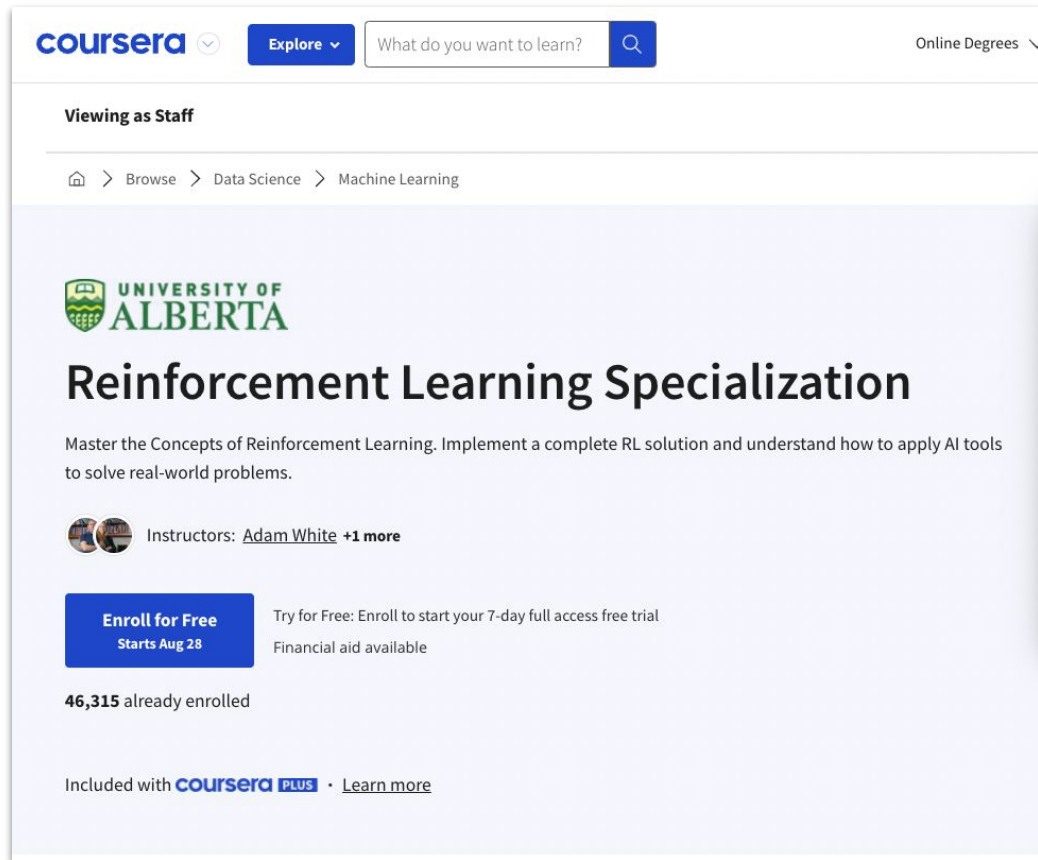
Assessment	Weight
Practice quizzes (80% pass)	9 x 1% = 9%
Assessments (graded quizzes/notebooks on Coursera)	9 x 2.5% = 22.5%
Midterm 1 exam	20 %
Midterm 2 exam	20%
Final exam	30%

Total: **101.5%**.
You can **not-submit**
2 quizzes and
2 assessments.

Grades **will not** be rounded at the end, and **no more extra marks** will be given.
No exceptions.

Coursera

- Coursera will be essential to CMPUT 365
- You should have been added to a private session of the RL courses (we used your university's email)
 - If you don't have access you should let me know!
 - **IMPORTANT: If you don't use the private session you won't get credit for submitted work!**



The screenshot shows the Coursera website interface. At the top, there is a search bar with the text "What do you want to learn?" and a search icon. To the left of the search bar is the Coursera logo and an "Explore" button. To the right is a dropdown menu for "Online Degrees". Below the search bar, it says "Viewing as Staff". A breadcrumb trail shows "Home > Browse > Data Science > Machine Learning". The main content area features the University of Alberta logo and the course title "Reinforcement Learning Specialization". Below the title is a description: "Master the Concepts of Reinforcement Learning. Implement a complete RL solution and understand how to apply AI tools to solve real-world problems." There are two instructor profile pictures and the text "Instructors: Adam White +1 more". A blue button says "Enroll for Free Starts Aug 28". To the right of the button, it says "Try for Free: Enroll to start your 7-day full access free trial" and "Financial aid available". Below the button, it says "46,315 already enrolled". At the bottom, it says "Included with Coursera PLUS · Learn more".

Fundamentals of Reinforcement Learning

Course Material

- Week 1
- Week 2
- Week 3
- Week 4

Grades

Notes

Discussion Forums

Messages

Live Events

Classmates

Course Manager Staff & Mentors Only

Welcome to the Course!

An Introduction to Sequential Decision-Making

All videos completed All readings completed All graded assessments completed

For the first week of this course, you will learn how to understand the exploration-exploitation trade-off in sequential decision-making, implement incremental algorithms for estimating action-values, and compare the strengths and weaknesses to...

Show Learning Objectives

The K-Armed Bandit Problem

- Module 1 Learning Objectives
Reading • 10 min
- Weekly Reading
Reading • 30 min
- Let's play a game!
Ungraded Plugin • 15 min
- Sequential Decision Making with Evaluative Feedback
Video • 5 min
- Compare bandits to supervised learning
Discussion Prompt • 10 min

What to Learn? Estimating Action Values

- Learning Action Values
Video • 4 min
- What's underneath?
Ungraded Plugin • 15 min
- Estimating Action Values Incrementally
Video • 5 min



Academic integrity

- [Code of Student Behaviour](#)
- [Student Conduct Policy](#)
- [Academic Integrity website](#)

- **Appropriate collaboration:** You are allowed to discuss the quizzes and assignments with your classmates. Note, however, that you are not allowed to exchange any written text, code, or to give and/or receive detailed step-by-step instructions on how to solve the proposed problems.

- **Cell phones:** Cell phones are to be turned off during lectures, labs and seminars.

- **Recording and/or Distribution of Course Materials:** Audio or video recording, digital or otherwise, by students is allowed only with my prior written consent as a part of an approved accommodation plan.

Academic integrity – **Expectations for AI use**

The primary goal of this course is to foster *individual* critical, creative thinking, and problem-solving skills related to reinforcement learning. Thus, in order to achieve such learning outcomes, you can submit each practice quiz and graded assignment multiple times, which allows for many learning opportunities.

The use of advanced AI-tools based on large-language models such as ChatGPT is **strictly prohibited** for all quizzes and graded assignments. The only exception is their use for Python-related queries (but the use of such tools to help with the programming assignments themselves is still strictly prohibited).

As stated in the university's [AI-Squared - Artificial Intelligence and Academic Integrity](#) webpage, *“learning is not only about the product; learning is also about the process of acquiring new knowledge or learning ways to think and reason.”*

Schedule

- The course will be structured in “weeks”. **Not every week starts on Monday**
- We have 12 weeks of content classes and we’ll cover 13 weeks of the MOOC
 - This corresponds to 9 chapters of the textbook

Schedule

- The course will be structured in “weeks”. **Not every week starts on Monday**
- We have 12 weeks of content classes and we’ll cover 13 weeks of the MOOC
 - This corresponds to 9 chapters of the textbook
- My overall (and tentative) plan for each one of the 3 days of the course-week:
 - 1st day: Non-comprehensive summary of the topic of the week
 - 2nd day: Non-comprehensive summary of the topic of the week + Questions + Exercises
 - 3rd day: Additional exercises, some in class activities

Schedule

- The course will be structured in “weeks”. **Not every week starts on Monday**
- We have 12 weeks of content classes and we’ll cover 13 weeks of the MOOC
 - This corresponds to 9 chapters of the textbook
- My overall (and tentative) plan for each one of the 3 days of the course-week:
 - 1st day: Non-comprehensive summary of the topic of the week
 - 2nd day: Non-comprehensive summary of the topic of the week + Questions + Exercises
 - 3rd day: Additional exercises, some in class activities
- A practice quiz is due in the 1st day of almost every course-week
- A graded assignment is due in the 2nd day of almost every course-week
- The deadline for submitting assignments and quizzes is 23:59:59

Schedule

Course Schedule & Assigned Readings				
Week	Date	Topic	Deadlines (all due at 23:59:59)	Readings
1	Wed, Sep 4	Course overview Discussion about what is reinforcement learning		
1	Fri, Sep 6	Background review: Probability, statistics, linear algebra, and calculus		
2	Mon, Sep 9	Fundamentals of RL: An introduction to sequential decision-making	Practice quiz (Sequential decision-making)	Chapter 2, up to §2.7 (pp. 25-36), and §2.10 (pp. 42-44)
2	Wed, Sep 11	Fundamentals of RL: An introduction to sequential decision-making	Program. assignment (Bandits & exploration / exploitation)	
3	Fri, Sep 13	Fundamentals of RL: Markov decision processes (MDPs)	Practice quiz (MDPs)	Chapter 3, up to §3.3 (pp. 47-56)
3	Mon, Sep 16	Fundamentals of RL: Markov decision processes (MDPs)		
4	Wed, Sep 18	Fundamentals of RL: Value functions & Bellman equations	Practice quiz (Value functions & Bellman equations)	Chapter 3, §3.5-§3.8 (pp. 58-69)
4	Fri, Sep 20	Fundamentals of RL: Value functions & Bellman equations	Graded quiz (Value functions & Bellman equations)	
4	Mon, Sep 23	Fundamentals of RL: Value functions & Bellman equations		
5	Wed, Sep 25	Fundamentals of RL: Dynamic programming	Practice quiz (Dynamic programming)	Chapter 4, §4.1-§4.4 (pp. 73-84); §4.6-§4.7 (pp. 86-89)
5	Fri, Sep 27	Fundamentals of RL: Dynamic programming	Program. Assignment (Optimal policies with dynamic programming)	
	Mon, Sep 30	National Day for Truth and Reconciliation		
5	Wed, Oct 2	Fundamentals of RL: Dynamic programming		

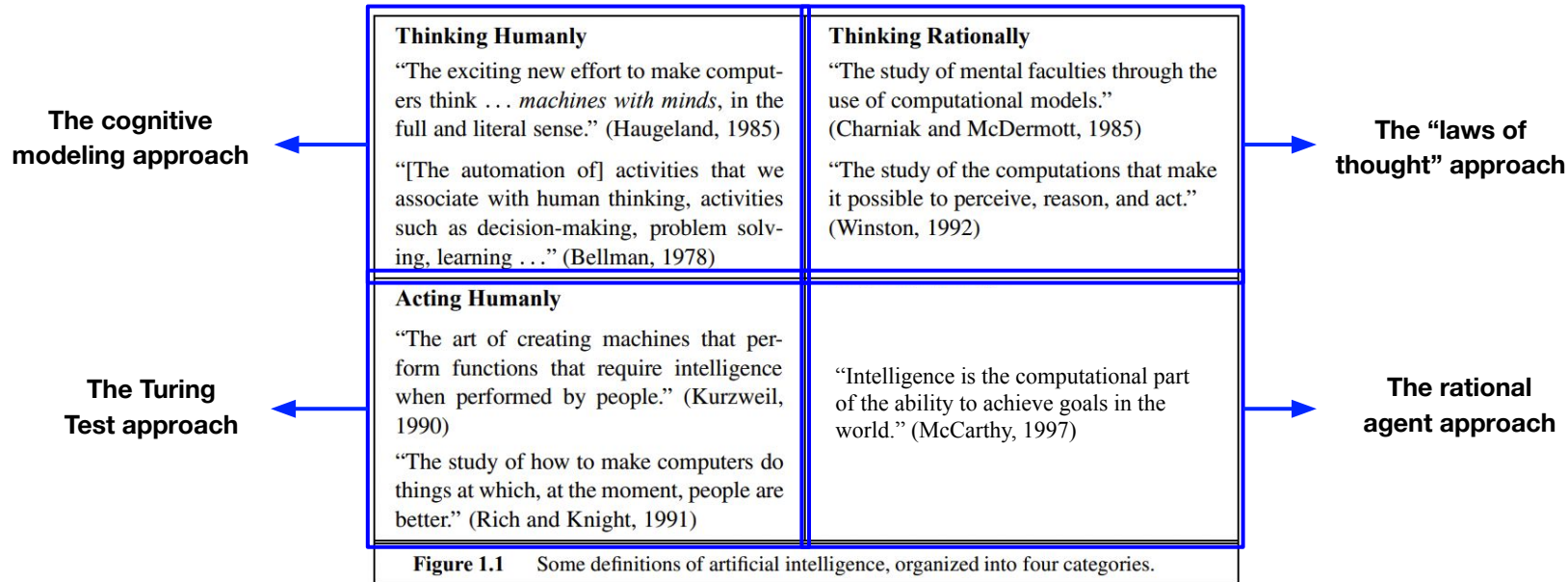
Syllabus [[eClass](#), [Slack](#), [website](#), [Google Drive](#)]



What is reinforcement learning?

Artificial intelligence

“AI is the ability of machines to perform tasks that are typically associated with human intelligence, such as learning and problem-solving.” –Wikipedia



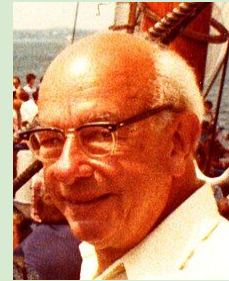
(Russell & Norvig; 2010)

Artificial intelligence

“AI is the ability of machines to perform tasks that are typically associated with human intelligence, such as learning and problem-solving.” –Wikipedia

The less a science has advanced, the more its terminology tends to rest on an uncritical assumption of mutual understanding.

– W. V. Quine



A diagram consisting of two nested rounded rectangles. The outer rectangle is light red and contains the text "Artificial intelligence". The inner rectangle is light green and contains the text "Machine learning". This visualizes machine learning as a subset of artificial intelligence.

Artificial intelligence

Machine learning

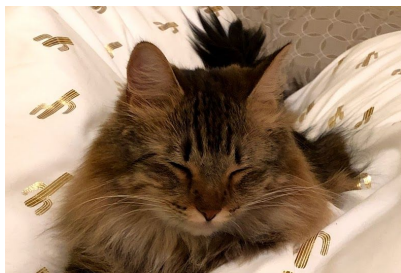
Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)



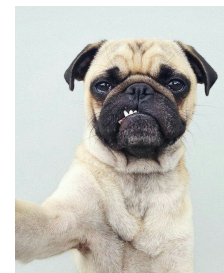
Cat



Cat



Not cat



Cat or not cat?

Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)
- “*Unsupervised learning* is typically about finding structure hidden in collections of unlabeled data” (Sutton & Barto; 2018)



Machine learning

Machine learning is a subfield of AI in which the system's desired behavior is not explicitly programmed, instead it is *learned* from data

- “*Supervised learning* is learning from a training set of labeled examples provided by a knowledgeable external supervisor” (Sutton & Barto; 2018)
- “*Unsupervised learning* is typically about finding structure hidden in collections of unlabeled data” (Sutton & Barto; 2018)



... and *reinforcement learning*!

Reinforcement learning

Artificial intelligence

Machine learning

Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)



Artificial intelligence

Machine learning

Reinforcement learning

Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)

- The idea of learning by interacting with our environment is very natural
- It is based on the idea of a learning system that wants something, and that adapts its behavior to get that



Artificial intelligence

Machine learning

Reinforcement learning

Reinforcement learning

Reinforcement learning is a computational approach to learning from interaction to maximize a numerical reward signal (Sutton & Barto; 2018)

- The idea of learning by interacting with our environment is very natural
- It is based on the idea of a learning system that wants something, and that adapts its behavior to get that



Some features are unique to reinforcement learning:

- Trial-and-error
- The trade-off between exploration and exploitation
- The delayed credit assignment / delayed reward problem

Reinforcement learning

Reinforcement learning is a computational paradigm for learning from interaction to maximize a numerical reward signal (Sutton & Barto, 2018)

- The idea of learning by interacting with an environment is very natural
- It is based on the idea of a human child that wants something, and that needs to learn how to get that

Some features are unique to reinforcement learning:

- Trial-and-error learning
- The trade-off between exploration and exploitation
- The delayed credit assignment / delayed reward problem

Artificial intelligence

Machine learning

Reinforcement learning

Problem or solution?



RL is now commonly deployed in the real-world

- **Recommendation systems**
 - Ads, news articles, videos, etc
- **General game playing**
 - Go, Chess, Shogi, Atari 2600, Starcraft, Minecraft, Gran Turismo
- **Industrial automation**
 - Cooling commercial buildings
 - Inventory management
 - Gas turbine optimization
 - Optimizing combustion in coal-fired power plants
- **Algorithms**
 - Video compression on YouTube
 - Faster matrix multiplication
 - Faster sorting algorithms
- **Control / Robotics**
 - Navigating stratospheric balloons
 - Plasm control for nuclear fusion
- **And more (see Csaba's [slides](#))**
 - COVID-19 border testing
 - Conversational agents
 - ...

On intelligence, AGI, etc etc...

- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence

On intelligence, AGI, etc etc...

- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence
- RL was originally developed to understand intelligence/the brain
 - We should develop a critical view around these topics, and an ability to recognize hype / PR pieces



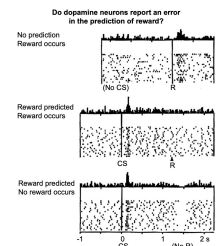
A. Barto (2024)

On intelligence, AGI, etc etc...

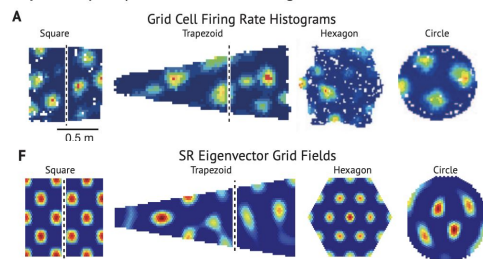
- People in the field have different, non-competing, perspectives and motivations
 - Some study RL to learn about / develop tools for solving sequential decision-making problems
 - Some look at RL as a computational model of intelligence
- RL was originally developed to understand intelligence/the brain
 - We should develop a critical view around these topics, and an ability to recognize hype / PR pieces
- Both perspectives are valid and both had had successes in the past **But they are different!!**



A. Barto (2024)



(Schultz, Dayan, & Montague; 1997)



(Stachenfeld, Botvinick, & Gershman; 2017)

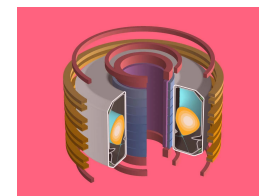
(Silver et al.; 2016)



(Degraeve et al.; 2022)



(Bellemare et al.; 2020)



Next class

- What **I** plan to do: A reminder about the required theoretical background
 - Probability (e.g., expectations of random variables, conditional expectations)
 - Calculus (e.g., partial derivatives)
 - Linear algebra (e.g., vectors and matrices)
 - I won't remind / teach you Python.
- What I recommend **YOU** to do for next class:
 - Make sure you have access to Coursera, eClass, and Slack
 - Brush up whatever you feel you are rusty on in terms of background
 - Read Chapter 1 of the textbook (not mandatory)
 - Start “Fundamentals of RL: An introduction to sequential decision-making” on Coursera (Week 1)