# DURHAM COLLEGE
## SUCCESS MATTERS

# School of Business, IT and Management

# Reinforcement Learning

Insert Academic Year

| School-Program | Year | Semester |
|---|---|---|
| Honours Bachelor of Artificial Intelligence | 4 | 8 |

| | |
|---|---|
| **Course Code:** Click or tap here to enter text. | **Course Equiv. Code(s)**: Click or tap here to enter text. |
| **Total Course Hours**: 42<br>**Lecture Hours**: 3<br>**Lab/Tutorial Hours**: N/A | **Credit Value**: 3 |
| **Prerequisite**: Introduction to Machine Learning, Introduction to Artificial Neural Networks | |
| **Corequisite**: N/A | |

| Laptop Course:<br>Yes ☒ | No ☐ | |
|---|---|---|
| **Delivery Mode(s):**<br>In class ☒ | **Online ☐** | **Hybrid ☐** |

| Authorized by (Dean): | Date: |
|---|---|

| |
|---|
| **Prepared by**: Dr. Ehsan Amjadian, PhD in Cognitive Science, Natural Language Processing, Deep Learning |
| **Qualified to teach**: PhD in Cognitive Science, Natural Language Processing, Deep Learning, Machine Learning, Computer Science |

# Course Description:

This course is an introduction to the field of reinforcement learning (RL). RL is primarily concerned with understanding how a software agent learns to behave in an environment to maximize the reward. Students explore methods, algorithms, and theoretical and practical aspects of RL, including exploration and generalization, the definition of state space, action space, dynamics, rewards, on-policy and off-policy learning, value iteration and policy iteration for RL.

# Course Learning Outcomes:

| Course Specific Learning Outcomes (CLO) |
| --- |
| Students receiving a credit for this course will have demonstrated their ability to: |
| 1. Explain the functions of the terms and elements of reinforcement learning formalisms, notations, and equations. |
| 2.  Explain the circumstances when reinforcement learning is applicable. |
| 3. Discuss the advantages and disadvantages of reinforcement learning by contrasting to other areas in machine learning. |
| 4. Code reinforcement learning algorithms to automatically achieve a goal having learnt from rewards and the absence of rewards. |
| 5. Analyze the architecture of deep reinforcement learning models to determine modifications required to improve performance on different tasks. |
| 6. Analyze the characteristics of various reinforcement learning paradigms to determine the appropriate application. |

Students will be notified in writing of changes that involve the addition or deletion of learning outcomes or evaluations, prior to changes being implemented, as specified in the Course Outline Policy and Procedure at Durham College.

# Undergraduate Degree Standards:

This course will contribute to the achievement of the following Undergraduate Degree Standards, as outlined in the Ontario Qualifications Framework.

| Undergraduate Degree Standards | Course Learning Outcome Reference |
| --- | --- |
| Depth and Breadth of Knowledge | CLO6 |
| Conceptual & Methodological Awareness/Research and Scholarship | CLO3 |

| | |
|---|---|
| Communication Skills | CLO1 |
| Application of Knowledge | CLO4 |
| Professional Capacity/Autonomy | CLO5 |
| Awareness of Limits of Knowledge | CLO2 |

## Program Learning Outcomes:

This course will contribute to the achievement of the following Program Learning Outcomes:

| | |
|---|---|
| 1. | Evaluate data requirements and technical approaches for building an AI solution by using systems analysis to align to business and client concerns. |
| 3. | Build machine learning models by evaluating input data and identifying features that meet the needs of the project. |

## Evaluation Criteria:

| Evaluation | Course Learning Outcomes | Weighting |
|---|---|---|
| Quiz 1 | CLO1, CLO2, CLO6 | 15% |
| Assignment 1 | CLO4, CLO5, CLO6 | 10% |
| Quiz 2 | CLO1, CLO3, CLO6 | 15% |
| Assignment 2 | CLO4, CLO5, CLO6 | 10% |
| Quiz 3 | CLO1, CLO6 | 15% |
| Quiz 4 | CLO1, CLO5, CLO6 | 15% |
| Term Project | CLO1, CLO2, CLO3, CLO4, CLO5, CLO6 | 20% |
| **Total** | | **100%** |

Students will be notified in writing of changes that involve the addition or deletion of learning outcomes or evaluations, prior to changes being implemented, as specified in the Course Outline Policy and Procedure at Durham College.

## Required Text(s) and Supplies:

1. RT1: Sutton, R. S., Barto, A. G. (2018). [Reinforcement Learning: An Introduction](#). The MIT Press.

2. RT2: Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M.A. (2013). [Playing Atari with Deep Reinforcement Learning](#). ArXiv, abs/1312.5602.

3. RT3: Hasselt, H.V., Guez, A., & Silver, D. (2015). [Deep Reinforcement Learning with Double Q-learning](#). AAAI.

4. [RT4: Adam Paszke Reinforcement Learning DQN (PyTorch) Tutorial.](#)

## Classroom Equipment and Requirements:

1. Students need access to cloud computing (Amazon AWS, Microsoft Azure, or Google Cloud Platform) in order to implement these techniques most efficiently.

2. Computers to access the cloud resources mentioned above.

## Recommended Resources (purchase is optional):

1. Lapan, M. (2018). [Deep Reinforcement Learning Hands-On](#). Birmingham, UK: Packt Publishing. ISBN: 978-1-78883-424-7

2. RT4: Fujimoto, S., Hoof, H.V., & Meger, D. (2018). [Addressing Function Approximation Error in Actor-Critic Methods](#). ICML.

**DURHAM COLLEGE**
SUCCESS MATTERS

| Learning Plan | | |
|---|---|---|
| **Week** | **Lecture** | **Hours:** |
| 1 | **Topics:**<br>● Course Syllabus & Outline<br>    ○ Course Materials<br>    ○ Topics per session<br>    ○ Evaluation<br>● What is Reinforcement Learning<br>    ○ Components of Reinforcement Learning<br>    ○ Examples<br>    ○ Early History<br>● Multi-armed Bandits Part 1<br><br>**Intended Learning Objectives:**<br>● Understand the RL problem<br>● Analyze the Multi-Armed Bandit problem | |
| | **Intended Learning Activities:**<br>● Lecture | |
| | **Resources and References:**<br>● RT1 Chapter 1 and 2 | |
| | **Evaluation and Weighting:**<br>N/A | |
| | **Course Learning Outcome Reference:**<br>CLO3 | |
| **Week** | **Lecture** | **Hours:** |
| 2 | **Topics:**<br>● Multi-armed Bandits Part 2<br>● Finite Markov Decision Processes Part 1<br><br>**Intended Learning Objectives:**<br>● Evaluate various multi-armed bandit problems and their corresponding solutions<br>● Understanding the components of a Markov Decision Process | |

| | | **Intended Learning Activities:**<br>● Lecture, in-class exercise |
|---|---|---|
| | | **Resources and References:**<br>● RT1 Chapter 2 and 3 |
| | | **Evaluation and Weighting:**<br>N/A |
| | | **Course Learning Outcome Reference:**<br>CLO1, CLO6 |
| **Week** | **Lecture** | **Hours:** |
| 3 | | **Topics:**<br>● Finite Markov Decision Processes Part 2<br>● Dynamic Programming<br><br>**Intended Learning Objectives:**<br>● Analyze Markov Decision Processes<br>● Evaluate Dynamic Programming Algorithms in computing optimal policies given a perfect model of the environment as a Markov decision process |
| | | **Intended Learning Activities:**<br>● Lecture, in-class exercise |
| | | **Resources and References:**<br>● RT1 Chapter 3 and 4 |
| | | **Evaluation and Weighting:**<br>N/A |
| | | **Course Learning Outcome Reference:**<br>CLO2 |
| **Week** | **Lecture** | **Hours:** |
| 4 | | **Topics:**<br>● Quiz 1<br>● Monte Carlo Methods<br><br>**Intended Learning Objectives:** |

| | |
|---|---|
| | ● Evaluate Monte Carlo Methods in finding optimal policies at the absence of complete knowledge of the environment |
| | **Intended Learning Activities:**<br>● Lecture, Quiz |
| | **Resources and References:**<br>● RT1 Chapter 5 |
| | **Evaluation and Weighting:**<br>Quiz 1 (15%) |
| | **Course Learning Outcome Reference:**<br>CLO1, CLO3 |
| **Week** | **Lecture**  **Hours:** |
| 5 | **Topics:**<br>● Temporal Difference Learning<br><br>**Intended Learning Objectives:**<br>● Evaluate TD Learning regarding how it approaches the prediction problem<br>● Compare and contrast TD against Monte Carlo Methods and Dynamic Programming |
| | **Intended Learning Activities:**<br>● Lecture, in-class exercise |
| | **Resources and References:**<br>● RT1 Chapter 6 |
| | **Evaluation and Weighting:**<br>N/A |
| | **Course Learning Outcome Reference:**<br>CLO6 |
| **Week** | **Lecture**  **Hours:** |
| 6 | **Topics:**<br>● n-step Bootstrapping<br>**Intended Learning Objectives:** |

| | |
|---|---|
| | ● Contrast one-step TD with n-step TD generalizing TD as well as Monte Carlso methods. |
| | **Intended Learning Activities:**<br>● Lecture, in-class exercise |
| | **Resources and References:**<br>● RT1 Chapter 6 |
| | **Evaluation and Weighting:**<br>Assignment 1 (10%) |
| | **Course Learning Outcome Reference:**<br>CLO1, CLO2, CLO3, CLO4 |
| **Week** | **Lecture**                                **Hours:** |
| 7 | **Topics:**<br>● Model-based and Model-free Reinforcement Learning<br><br>**Intended Learning Objectives:**<br>● Compare and contrast planning and learning in reinforcement learning<br>● Evaluate model-based and model-free approaches in computing a value function |
| | **Intended Learning Activities:**<br>● Lecture, in-class programming exercise |
| | **Resources and References:**<br>● RT1 Chapter 8 |
| | **Evaluation and Weighting:**<br>N/A |
| | **Course Learning Outcome Reference:**<br>CLO1, CLO6 |
| **Week** | **Lecture**                                **Hours:** |
| 8 | **Topics:**<br>● Quiz 2<br>● Approximating On-policy Prediction<br><br>**Intended Learning Objectives:** |

![Durham College logo - SUCCESS MATTERS]

| | |
|---|---|
| | ● Analyze function approximation in reinforcement learning to estimate a state-value function.<br>● Contrast function approximation methods in RL with tabular reinforcement learning |
| | **Intended Learning Activities:**<br>● Lecture, in-class programming exercise |
| | **Resources and References:**<br>● RT1 Chapter 9 |
| | **Evaluation and Weighting:**<br>Quiz 2 (15%) |
| | **Course Learning Outcome Reference:**<br>CLO3, CLO6 |
| **Week** | **Lecture**            **Hours:** |
| 9 | **Topics:**<br>● On-policy Control with Approximation<br><br>**Intended Learning Objectives:**<br>● Evaluating the control problem with a parametric approximation of the action-value function in on-policy RL |
| | **Intended Learning Activities:**<br>● Lecture |
| | **Resources and References:**<br>● RT1 Chapter 10 |
| | **Evaluation and Weighting:**<br>N/A |
| | **Course Learning Outcome Reference:**<br>CLO1, CLO6 |
| **Week** | **Lecture**            **Hours:** |
| 10 | **Topics:**<br>● Off-policy Methods with approximation<br>**Intended Learning Objectives:** |

| | | |
|---|---|---|
| | ● Contrast off-policy with function approximation with its on policy counterparts<br>● Analyze off-policy RL with approximation in terms of convergence and performance. | |
| | **Intended Learning Activities:**<br>● Lecture, in-class programming | |
| | **Resources and References:**<br>● RT1 Chapter 11 | |
| | **Evaluation and Weighting:**<br>Assignment 2 (10%) | |
| | **Course Learning Outcome Reference:**<br>CLO3, CLO4, CLO6 | |
| **Week** | **Lecture** | **Hours:** |
| 11 | **Topics:**<br>● Eligibility Traces | |
| | **Intended Learning Objectives:**<br>● Lecture, in-class exercise | |
| | **Intended Learning Activities:**<br>● Analyze the spectrum created by the eligibility traces in order to recognize the extreme poles and find the optimal methods<br>● Evaluate how eligibility traces combine Monte Carlo methods and TD methods | |
| | **Resources and References:**<br>● RT1 Chapter 12 | |
| | **Evaluation and Weighting:**<br>N/A | |
| | **Course Learning Outcome Reference:**<br>CLO6 | |
| **Week** | **Lecture** | **Hours:** |
| 12 | **Topics:**<br>● Quiz 3 | |

- Policy Gradient Methods

**Intended Learning Objectives:**
- Evaluate parametrized policies that directly select actions without consulting a value function.
- Analyze how might value functions interact with such parameterized policies.

**Intended Learning Activities:**
- Lecture

**Resources and References:**
- RT1 Chapter 13

**Evaluation and Weighting:**
Quiz 3 (15%)

**Course Learning Outcome Reference:**
CLO4, CLO6

| Week | Lecture | Hours: |
|------|---------|--------|
| 13 | **Topics:**<br>- Neural Actor-Critic<br>- TD-Gammon<br>- DQN<br>- Double Q-learning<br><br>**Intended Learning Objectives:**<br><br>**Intended Learning Activities:**<br>- Analyze TD-Gammon and modern neural Q-learning, comparing the two paradigms, advantages and disadvantages<br><br>**Resources and References:**<br>- RT1 Chapter 16<br>- RT2<br>- RT3<br><br>**Evaluation and Weighting:**<br>N/A | |

**DURHAM COLLEGE**
**SUCCESS MATTERS**

| | Course Learning Outcome Reference: CLO5, CLO6 | |
|---|---|---|
| **Week** | **Lecture** | **Hours:** |
| 14 | **Topics:**<br>● Term Project Due<br>● DQN with Pytorch<br>● Quiz 4<br><br>**Intended Learning Objectives:**<br>● Applying the PyTorch framework to implement deep reinforcement learning algorithms. | |
| | **Intended Learning Activities:**<br>● Facilitated group coding workshop<br>● Question & Answer | |
| | **Resources and References:**<br>● RT4 | |
| | **Evaluation and Weighting:**<br>Quiz 4 (15%)<br>Term Project (20%) | |
| | **Course Learning Outcome Reference:**<br>CLO4, CLO6 | |