

# RL Traffic Signal Control

---

David Sanwald

May 7, 2018

# Table of contents

1. Introduction

2. Reinforcement Learning

# Introduction

---

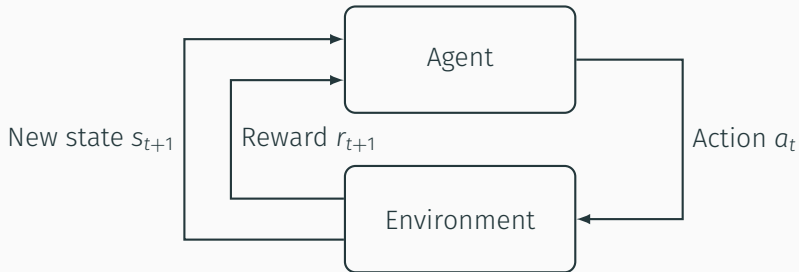
# Reinforcement Learning

---

# Machine Learning

- Supervised Learning
  - Classification
  - Regression
- Unsupervised Learning
  - Clustering
  - ...
- Reinforcement Learning

# Agent Environment



**Figure 1:** Agent environment interface

# Markov Decision Process

MARKOV DECISION PROCESS is defined by quatuple  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \mathcal{A} \rangle$

- $\mathcal{S}$ , a set of states
- $\mathcal{P}$ , a state transition matrix defining the probabilities of some possible next state  $s'$  given any state  $s$   
 $\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
- a reward function  $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- $\mathcal{A}$ , a set of actions

- specifies agent's behaviour
- mapping of state to action

$$\pi(s) = a$$

$$\mathbb{P}(a|s) = \pi(a|s)$$

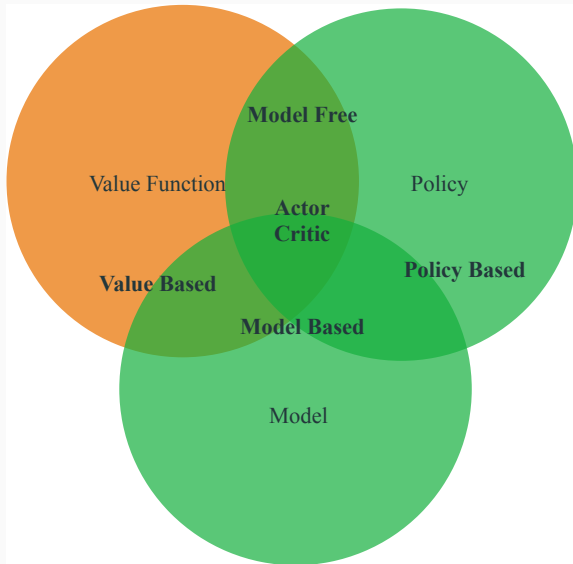


# Markov Property

- The future is conditionally independent of the past given the presence
- implies memorylessness

$$\mathbb{P}[S_{t+1}|S_1, \dots, S_t] = \mathbb{P}[S_{t+1}|S_t]$$

# Taxonomy of RL



# Value Function

Expected return

- from state  $s$  and action  $a$
- given policy  $\pi$

$$Q^\pi(s, a) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s, a]$$

- decomposable into

$$Q^\pi(s, a) = \mathbb{E}[r + \gamma Q^\pi(s', a') | s, a]$$

# Optimal Value Function

- optimal value function

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$$

- optimal policy

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

- decomposition into

$$Q^*(s, a) = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Off Policy learning

## Q-learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(S_t, A_t)]$$

On Policy learning

## Sarsa

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

Demo

# Function Approximation

Why Function Approximation?

- large state spaces
- slow learning
- need for generalization

# Deadly Triad

oho

Deadly Triad

oho